

SVM

Hendrik Santoso Sugiarto

IBDA2032 – *Artificial Intelligence*

Capaian Pembelajaran

- Linear SVM
- Nonlinear SVM
- Multi-class SVM
- Regresi SVM

Linear SVM

Support Vector Machine (SVM)

- Paradigma machine learning yang menggunakan perspektif pemisahan data (jarak antar data)
- Dapat digunakan untuk regresi, klasifikasi, unsupervised
- Cocok untuk klasifikasi data dengan banyak fitur (high dimension) dan jumlah data kecil / menengah (kurang efisien untuk data besar)
- Non-probabilistic: tidak memiliki nilai peluang

Formulasi

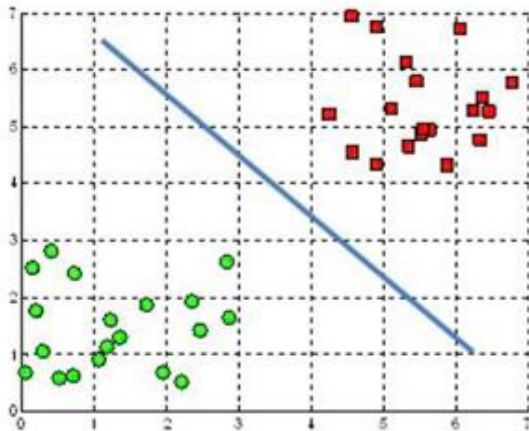
- Setting

- Data latih: $((\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(m)}, y^{(m)}))$
- Untuk klasifikasi biner: $y^{(i)} \in \{+1, -1\}$

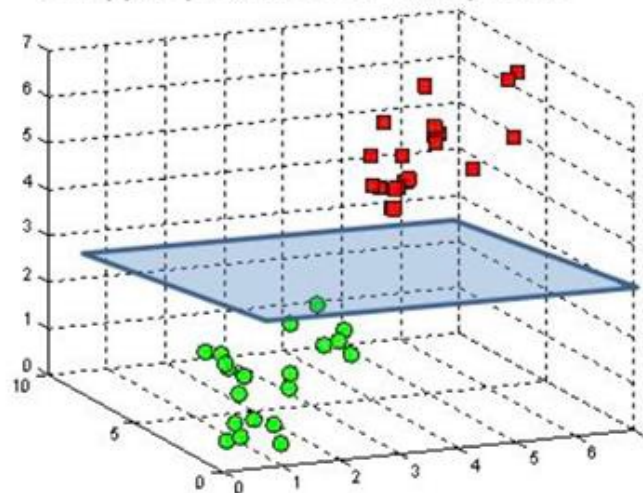
- Tujuan

- Menemukan hyperplane optimal yang dapat memisahkan data: $f(\mathbf{x}) = \text{sgn}(\mathbf{w}^T \mathbf{x} + b)$
- Hyperplane: subspace dengan dimensi $n - 1$ dari dimensi data

A hyperplane in \mathbb{R}^2 is a line



A hyperplane in \mathbb{R}^3 is a plane



$$0 = w_1 x_1 + w_2 x_2 + b$$
$$\frac{-w_1 x_1 - b}{w_2} = x_2$$

Uji Pemahaman

- Berapa dimensi hyperplane dari data dengan fitur berdimensi 1?

$$\left(\mathbf{x}^{(i)} = \{x_1^{(i)}\}\right) \quad 0$$

- Berapa dimensi hyperplane dari data dengan fitur berdimensi 2?

$$\left(\mathbf{x}^{(i)} = \{x_1^{(i)}, x_2^{(i)}\}\right) \quad 1$$

- Berapa dimensi hyperplane dari data dengan fitur berdimensi 3?

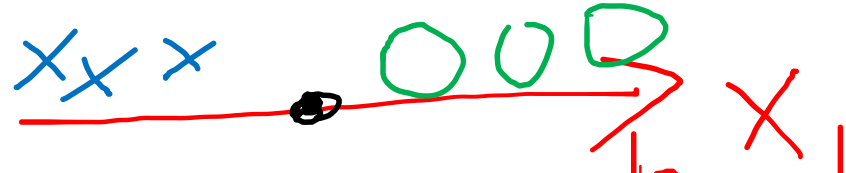
$$\left(\mathbf{x}^{(i)} = \{x_1^{(i)}, x_2^{(i)}, x_3^{(i)}\}\right) \quad 2$$

- Berapa dimensi hyperplane dari data dengan fitur berdimensi 5?

$$\left(\mathbf{x}^{(i)} = \{x_1^{(i)}, x_2^{(i)}, x_3^{(i)}, x_4^{(i)}, x_5^{(i)}\}\right) \quad 4$$

Uji Pemahaman

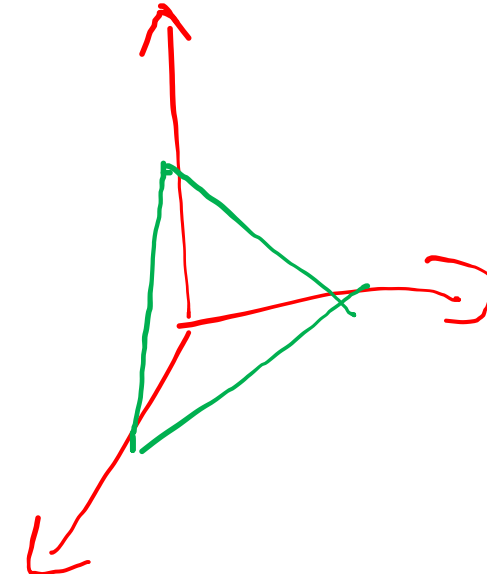
- Bagaimana bentuk persamaan hyperplane kasus sebelumnya?

$0 = w^T x + b$ → 1D = 

$w_1 x_1 + b = 0 \rightarrow x_1 = -\frac{b}{w_1}$

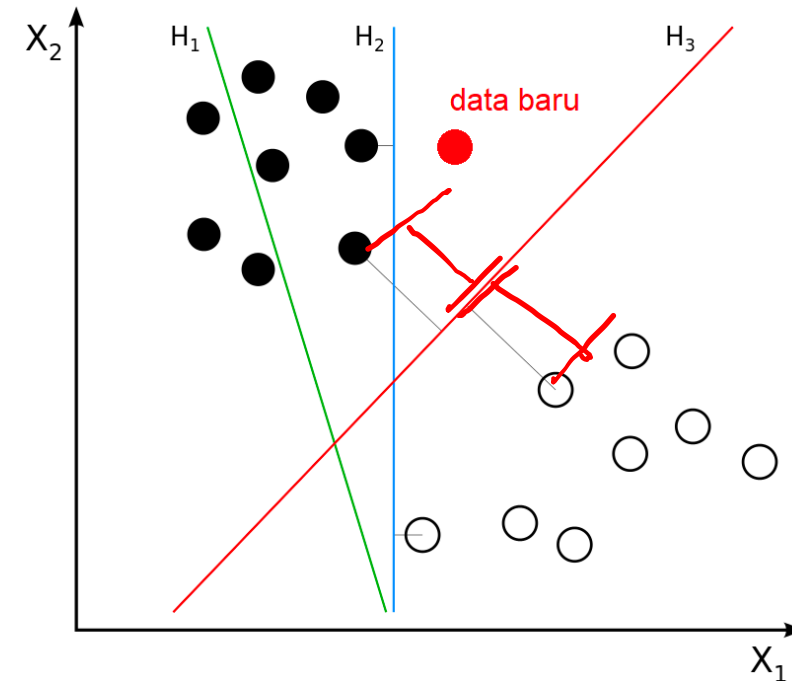
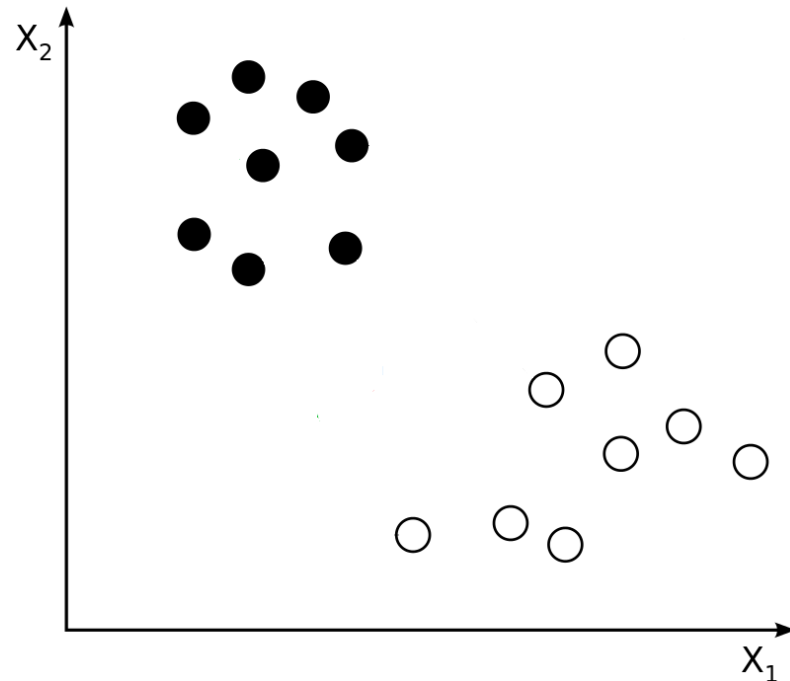
→ 2D = $x_1 = -\frac{w_2}{w_1} x_2 - \frac{b}{w_1}$

→ 3D = $-\frac{b}{w_1} - \frac{w_2}{w_1} x_2 - \frac{w_3}{w_1} x_3 = x_1$



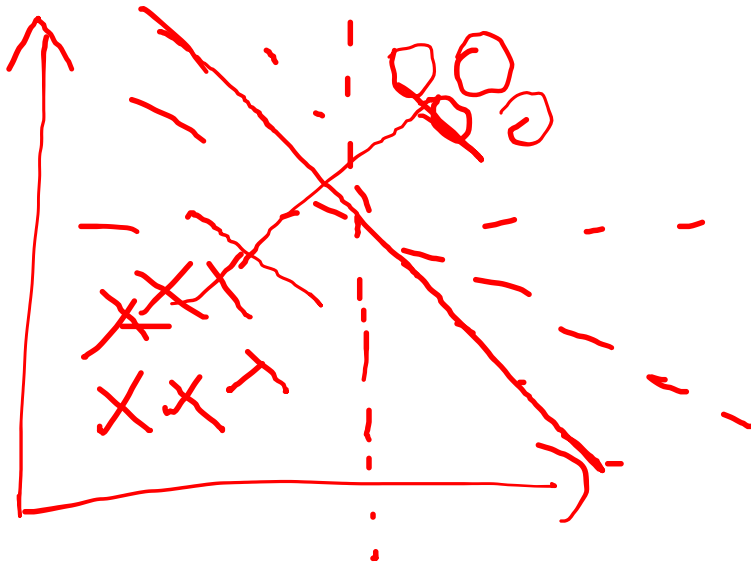
SVM Linear

- Bagaimana menarik garis lurus yang memisahkan 2 kelas berikut?
- H1 (hijau) → tidak dapat digunakan untuk klasifikasi
- H2 (biru) → mampu klasifikasi, garis pemisah dekat data
- H3 (merah) → mampu klasifikasi, garis pemisah jauh dari data



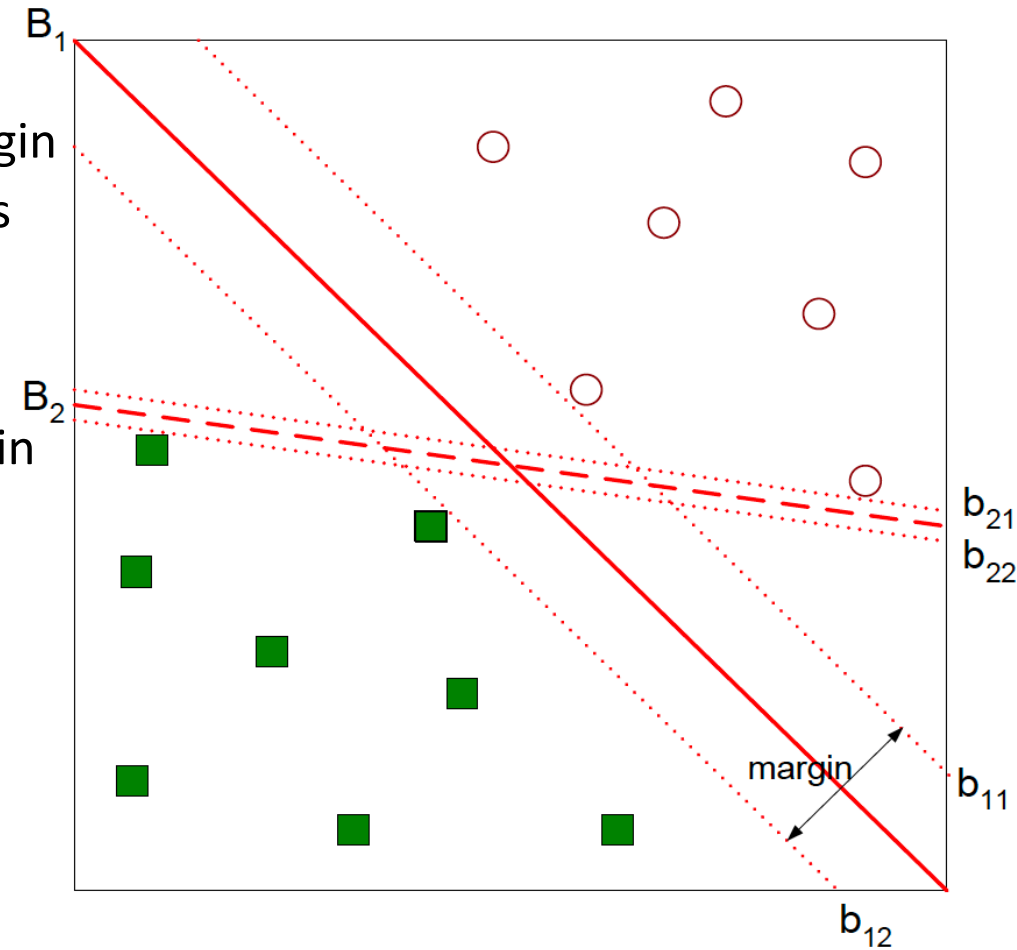
SVM Linear

- Manakah garis yang paling sesuai untuk memisahkan data? (bagaimana performa terhadap data baru?)
- Kelas yang berbeda dapat dengan mudah dipisahkan oleh lebih dari satu kemungkinan hyperplane (linearly separable)
- Metode SVM akan memilih garis yang berada sejauh mungkin dari data latih terdekat → large margin classifier



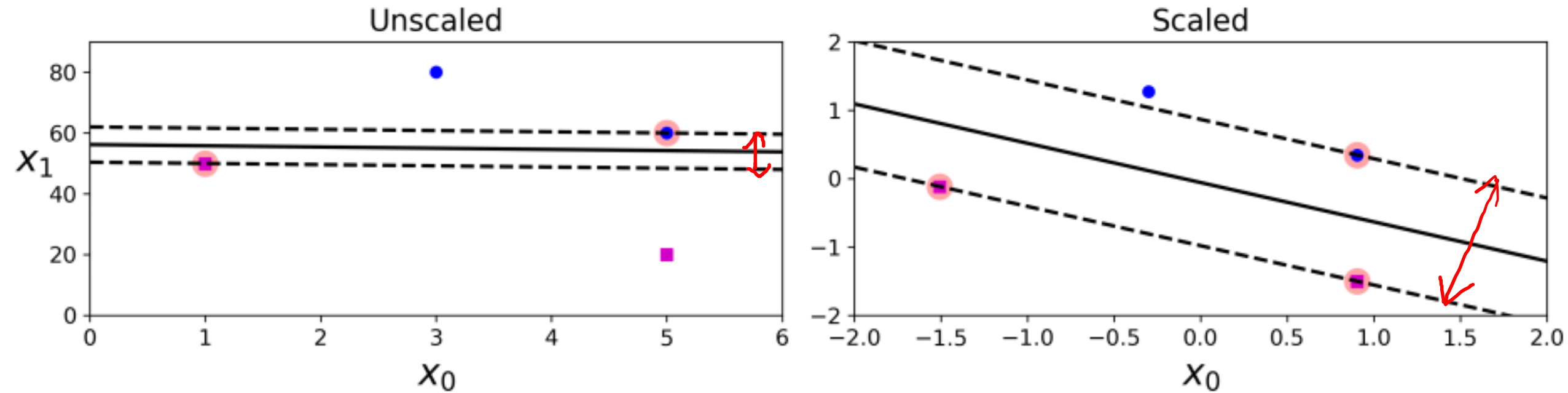
Intuisi: Margin Maksimum

- Intuisi sebuah margin:
 - Margin: lebar dari batas yang dapat diperbesar sebelum menyentuh sebuah data
- SVM:
 - Menemukan hyperplane yang dapat memaksimalkan margin
 - Batas keputusan harus sejauh mungkin dari data tiap kelas
- Support vector:
 - Data dari masing-masing kelas di tepi
 - Istilah: penambahan data latih lainnya berada diluar margin
 - → tidak mengubah Batasan keputusan
 - Batasan keputusan didukung (supported) oleh data di tepi



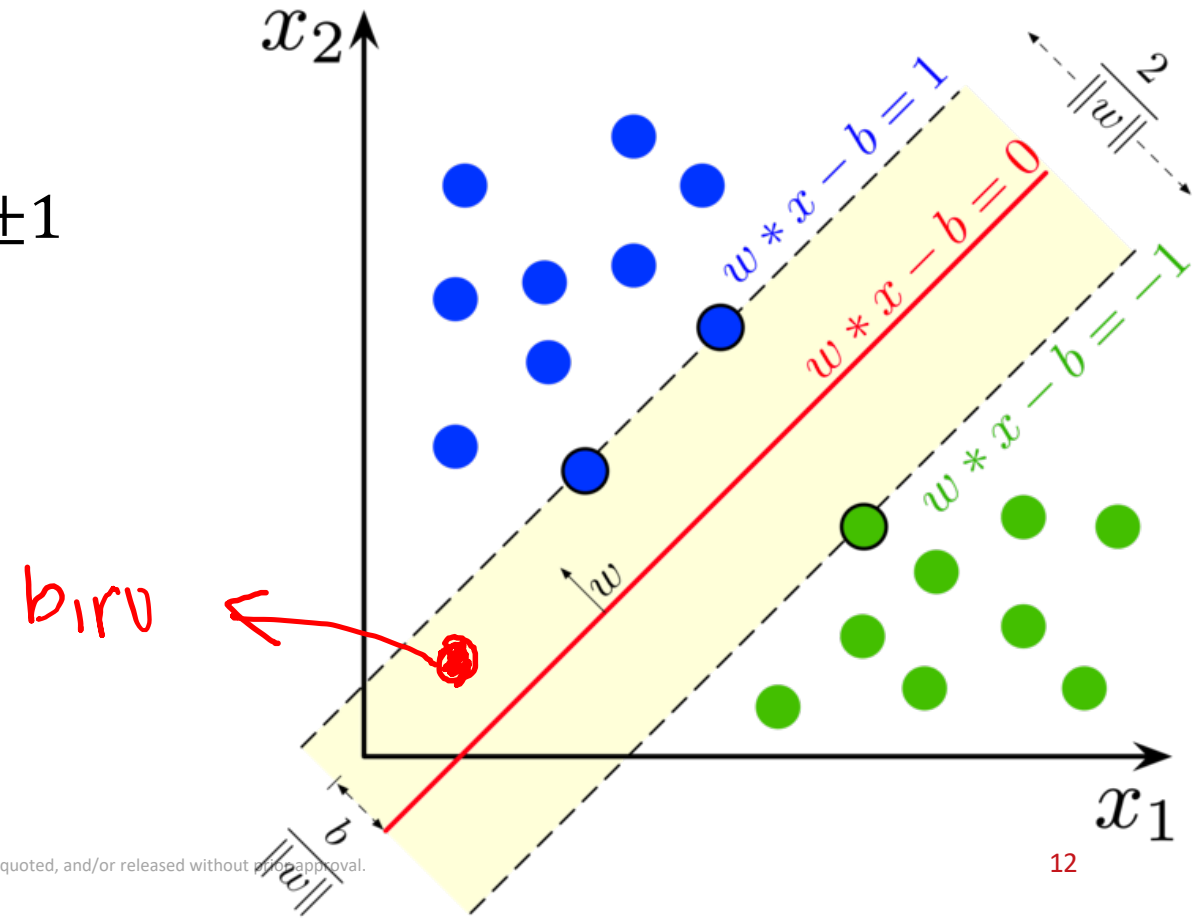
SVM sensitif terhadap scaling

- Unscaled: skala vertical jauh lebih besar dari horizontal
- Scaled: Batasan keputusan nampak lebih jelas



Prediksi SVM

- Klasifikasi SVM linear \rightarrow prediksi kelas dari data baru \mathbf{x} dengan menghitung fungsi keputusan h : $h = \mathbf{w}^T \mathbf{x} + b = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$
- Prediksi: $\hat{y} = \begin{cases} 0, & \text{jika } h < 0 \\ 1, & \text{jika } h \geq 0 \end{cases}$ -1ve
+1ve
- Batas keputusan: garis tegas dimana $h = 0$
- Batas margin: garis putus-putus dimana $h = \pm 1$



Optimisasi

lagrange multiplier

- Dataset: $\{x_1, \dots, x_n\}$, label: $y^{(i)} \in \{+1, -1\}$
- Batas keputusan harus sejauh mungkin dari data tiap kelas \rightarrow memaksimalkan margin:

$$m = \frac{2}{\|\mathbf{w}\|} \rightarrow \|\mathbf{w}\|^2 = \mathbf{w}^T \mathbf{w}$$

- Batasan keputusan harus meklasifikasi setiap data secara benar:

$$y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1, \forall_i$$

- A linear constrained convex quadratic optimization problem (quadratic programming):

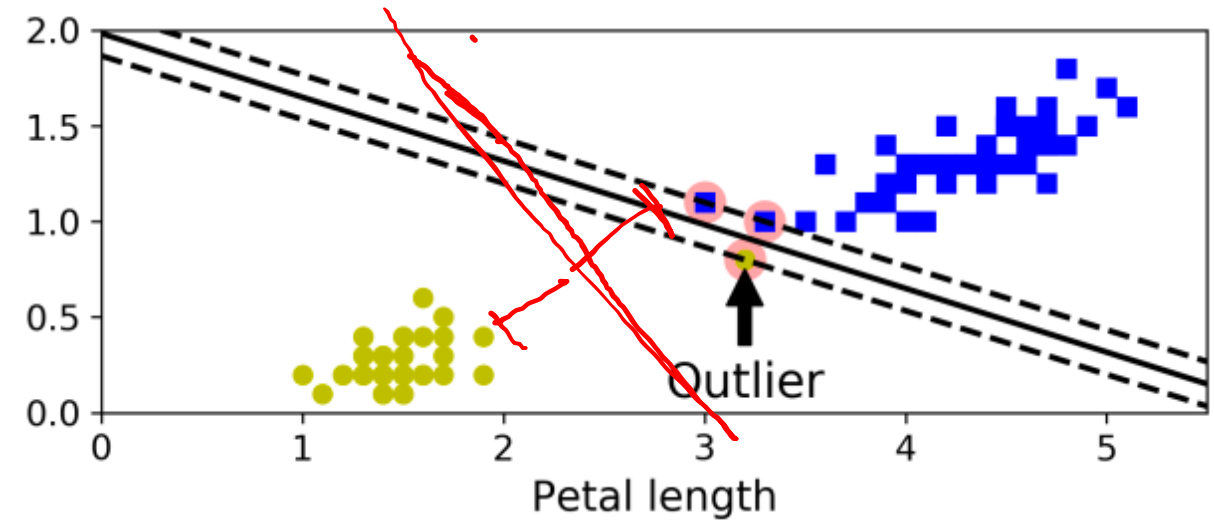
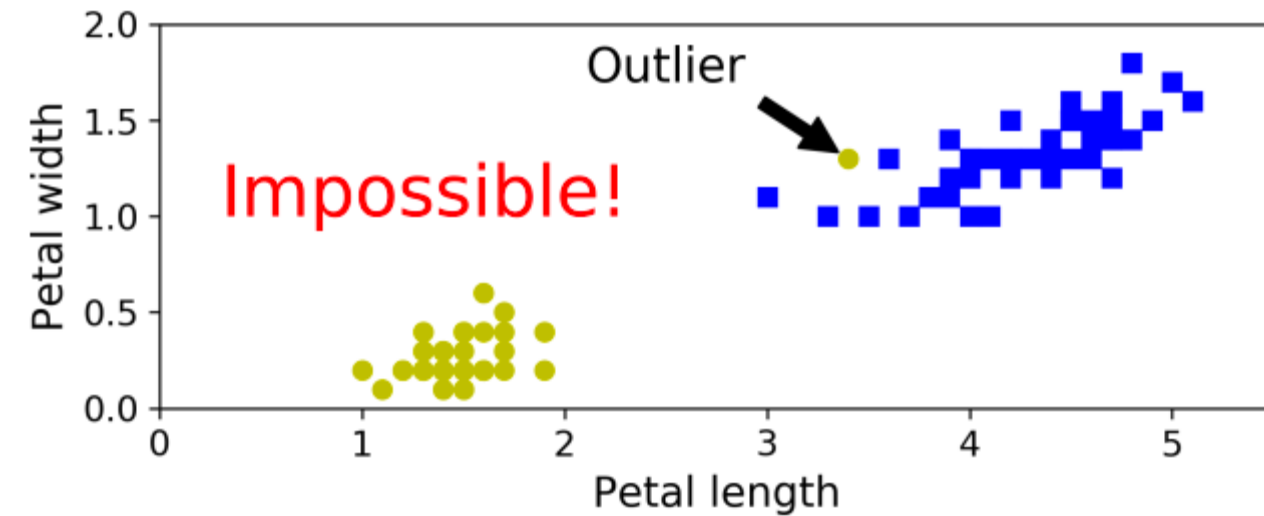
$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

subject to $\rightarrow \underline{y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1, \forall_i}$

constraint

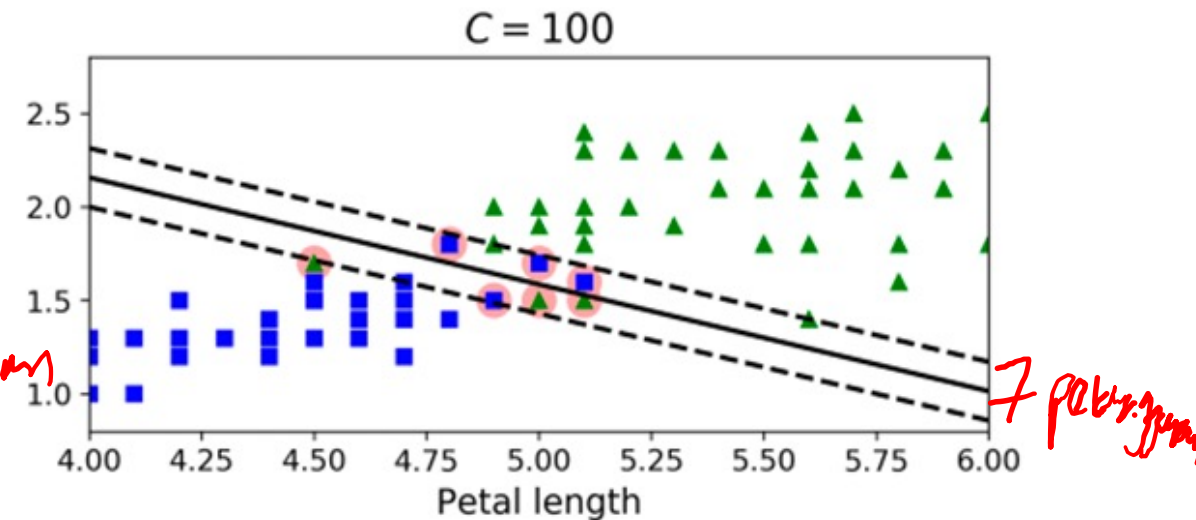
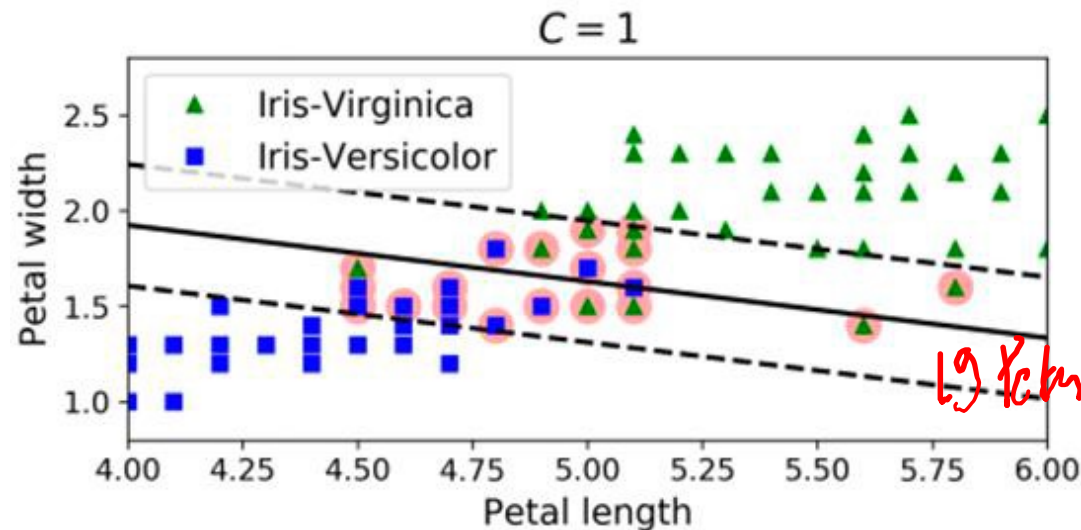
Hard Margin

- Hard margin → semua data harus di luar margin
- Masalah dengan hard margin:
 - Hanya bisa dilakukan jika data linearly separable
 - Sensitif terhadap pencilan: margin menjadi sempit



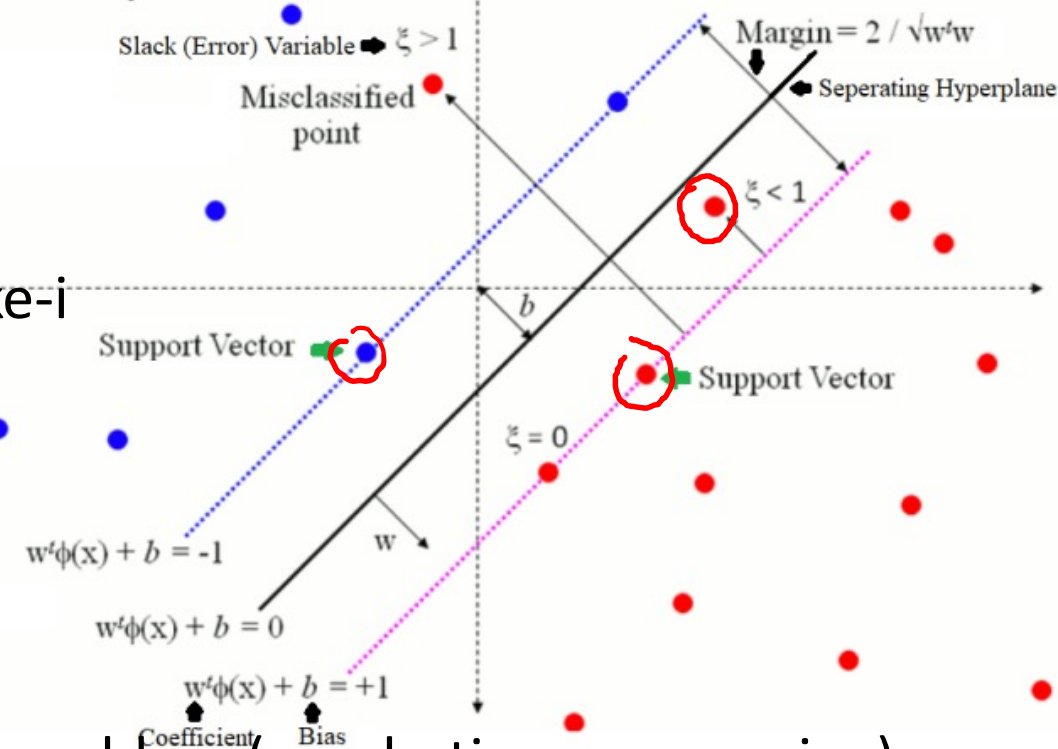
Soft Margin

- Soft margin merupakan kompromi antara:
 - margin selebar mungkin (margin optimization)
 - membatasi pelanggaran batas (margin violation)
- Margin violation: data yang melanggar batas margin
- Kompromi ini diatur oleh hyperparameter C (nilai kecil artinya margin besar)
- Kasus C mana yang cenderung lebih baik generalisasinya?



Optimisasi Soft Margin

- Menggunakan slack variables $\zeta^{(i)} \geq 0$ untuk data ke- i
- $\zeta^{(i)} \rightarrow$ data ke- i boleh melanggar batas sejauh apa
- 2 tujuan yang berseberangan:
 - $\zeta^{(i)}$ sekecil mungkin: menurunkan pelanggaran batas
 - $\frac{1}{2} ||\mathbf{w}'||^2$ sekecil mungkin: membesarkan margin
- A linear constrained convex quadratic optimization problem (quadratic programming):



$$\min_{\mathbf{w}, b, \zeta} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^m \zeta^{(i)}$$

margin \leftarrow $\frac{1}{2} \mathbf{w}^T \mathbf{w}$ \rightarrow pelanggaran

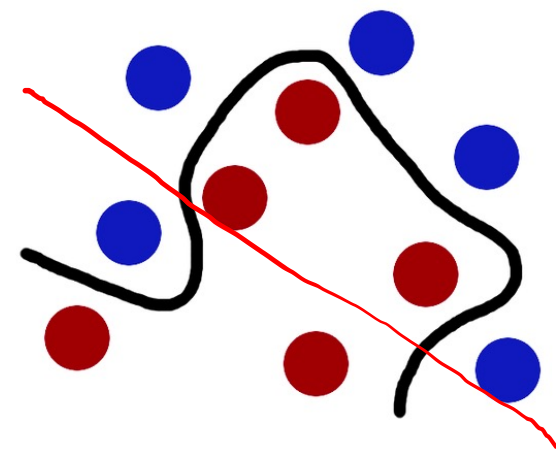
subject to $\rightarrow y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 - \zeta^{(i)}$ dan $\zeta^{(i)} \geq 0, \forall_i$

Hyperparameter

Non-Linear SVM

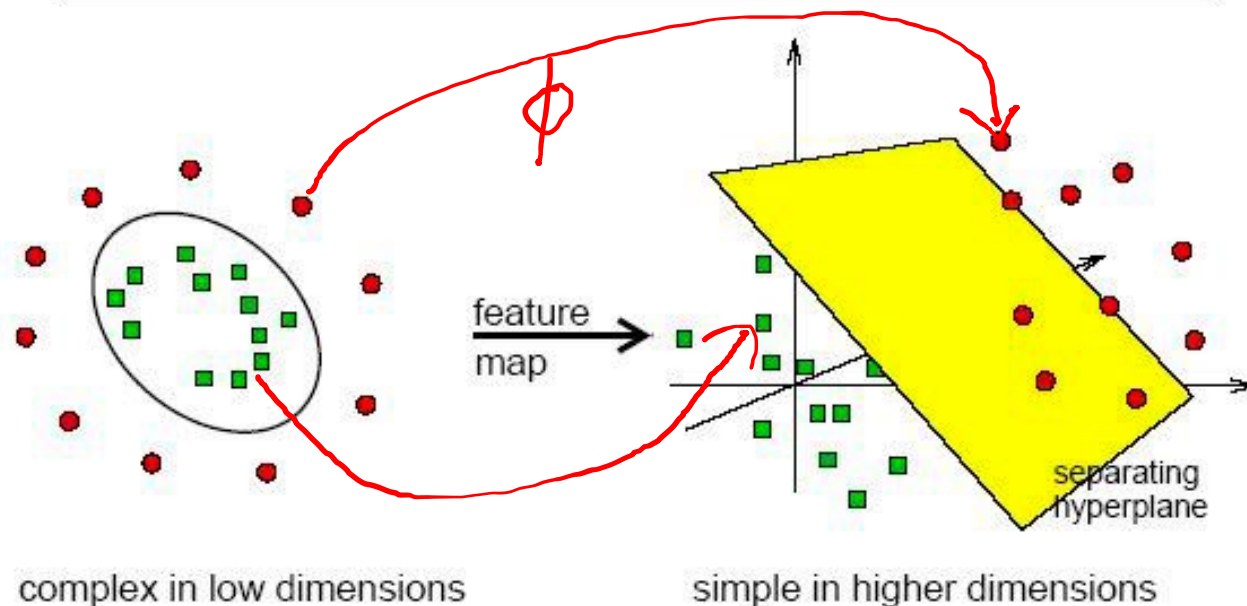
SVM NonLinear

$\phi(x)$



- Banyak data tidak terpisahkan secara linear
- Kernel trick: SVM linear untuk data yang non-linear
- Fungsi kernel → memetakan data non-linear ke dimensi yang lebih tinggi sehingga dapat terpisahkan secara linear
- Menambah fitur: misalnya dengan fitur polynomial lebih tinggi → mirip dengan regresi polynomial

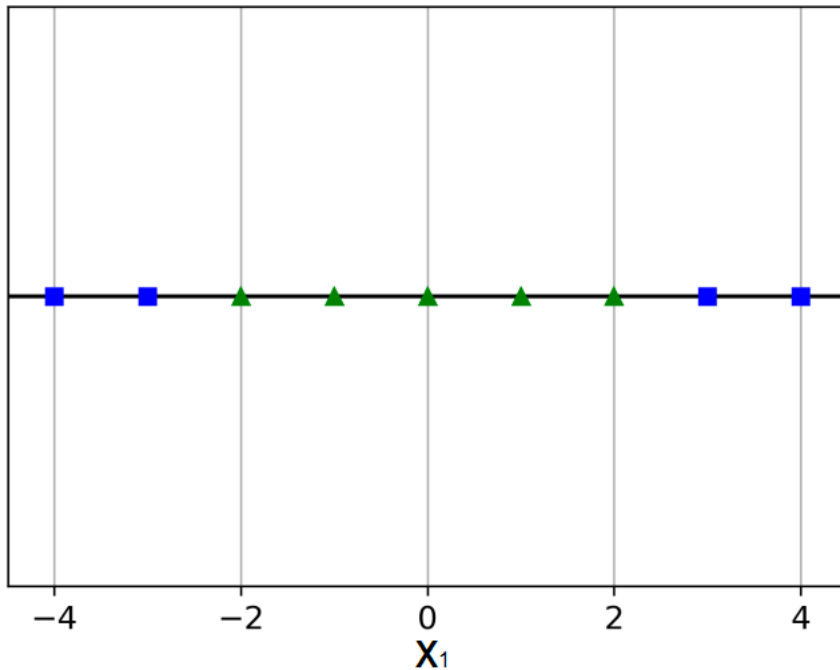
Separation may be easier in higher dimensions



Feature Mapping

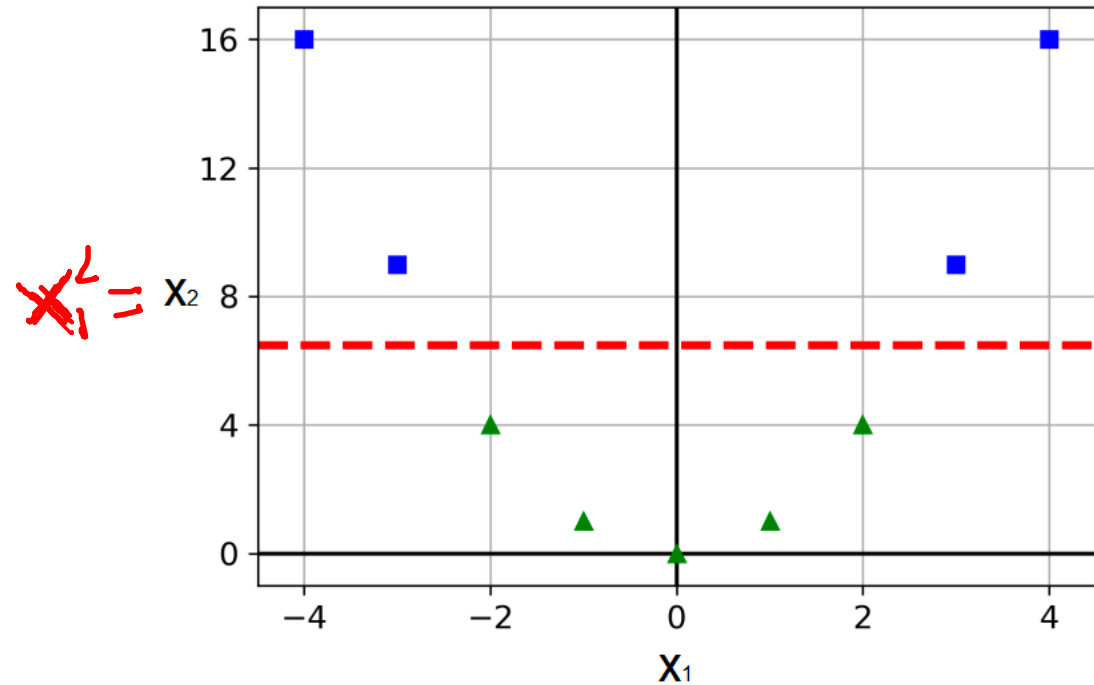
- Data dengan 1 fitur (x_1) \rightarrow tidak terpisahkan linear
- Tambahkan 1 fitur polynomial derajat 2 ($x_2 = x_1^2$) \rightarrow terpisahkan linear

1D

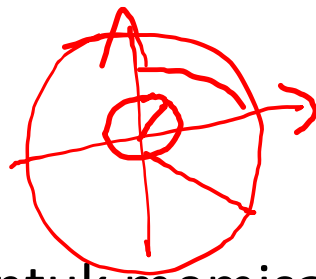


$\phi \rightarrow$

2D

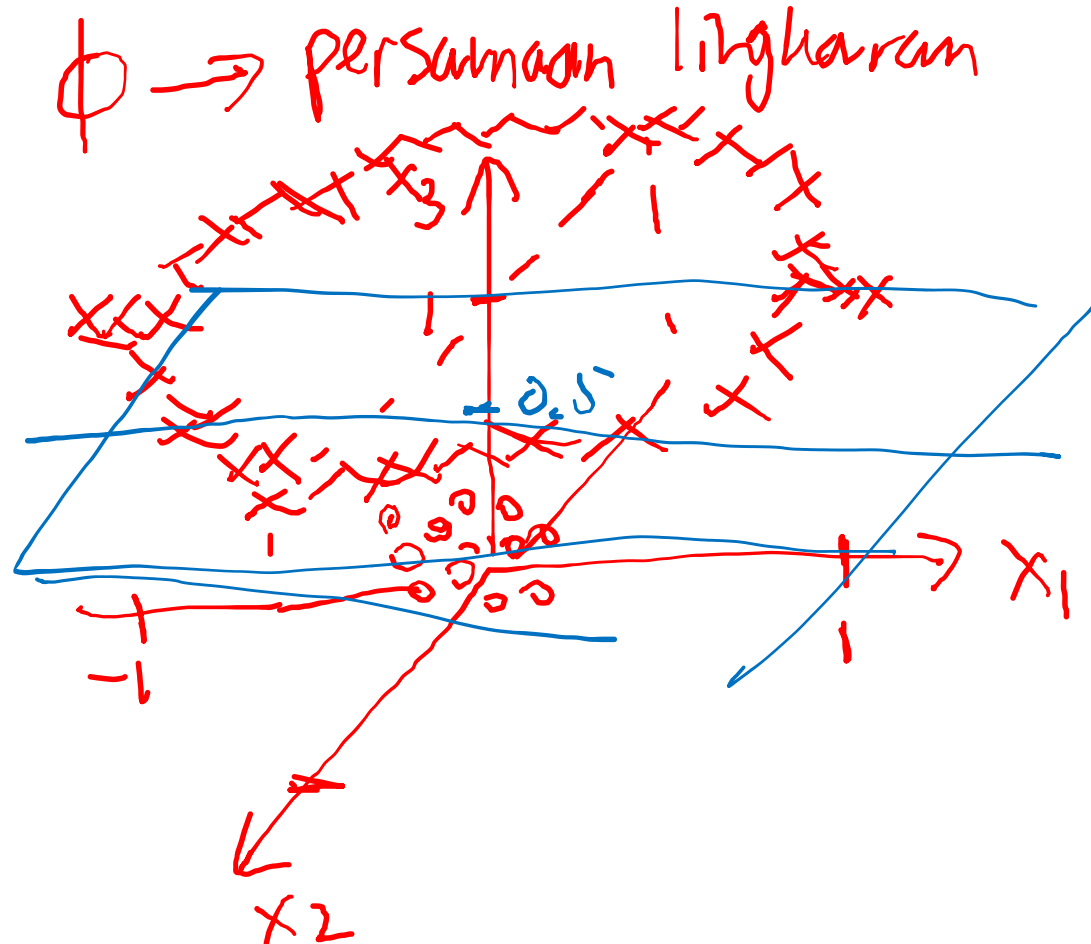
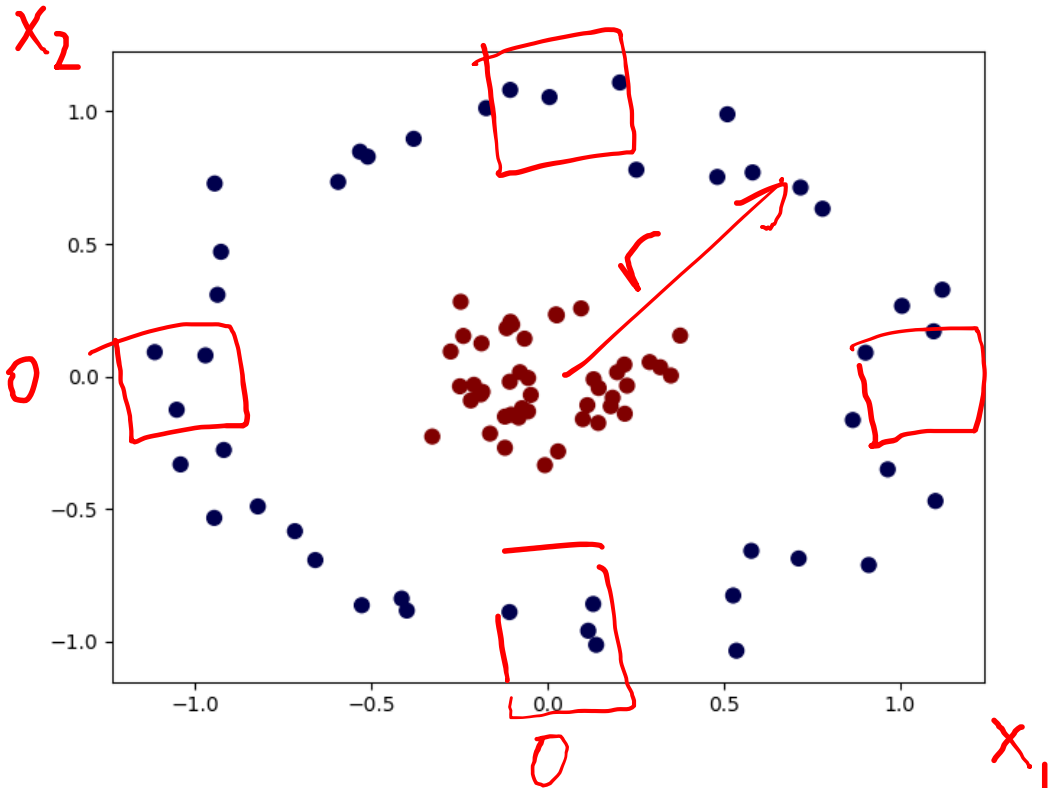


Uji Pemahaman



$$x_1^2 + x_2^2 = r^2 = x_3$$

- Gunakan polynomial mapping untuk memisahkan data berikut secara linear



Optimisasi SVM Nonlinear

- Setelah menggunakan feature mapping:

$$x \mapsto \Phi(x)$$

data
→ proyeksi data

- Temukan hyperplane pada feature space

$$f(x) = \mathbf{w} \cdot \Phi(x) + b$$

- Hampir sama dengan optimisasi SVM linear

$$\min_{\mathbf{w}, b, \zeta} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^m \zeta^{(i)}$$

$$\text{subject to } y^{(i)} (\mathbf{w}^T \phi(x^{(i)}) + b) \geq 1 - \zeta^{(i)} \text{ dan } \zeta^{(i)} \geq 0, \forall_i$$

Bagaimana cara memilih feature map

$$x \mapsto \phi(x)$$

- Fitur polynomial:
 - Jika derajat terlalu rendah → tidak dapat terpisahkan linear
 - Jika derajat terlalu tinggi → jumlah fitur besar → model terlalu lambat (butuh banyak memori dan sulit menyelesaikan quadratic programming)

• Contoh:

$x_1, x_2, x_3 \rightarrow \phi$ derajat 2

$$x \in R^3, \phi(x) \in R^{10}$$

$$\phi(x) = (1, x_1, x_2, x_3, x_1^2, x_2^2, x_3^2, x_1x_2, x_1x_3, x_2x_3)$$

Kernel Tricks

- SVM dengan fitur polynomial derajat tinggi dapat dipercepat dengan kernel trick
- Ide: ganti dot product dengan sebuah kernel

$$\kappa(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \langle \phi(\mathbf{x}^{(i)}), \phi(\mathbf{x}^{(j)}) \rangle$$

- Fungsi kernel:

- Linear kernel

$$\kappa(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \langle \mathbf{x}^{(i)}, \mathbf{x}^{(j)} \rangle = \mathbf{x}^{(i)T} \mathbf{x}^{(j)}$$

→ dot product

- Polynomial kernel

$$\kappa(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \langle \phi(\mathbf{x}^{(i)}), \phi(\mathbf{x}^{(j)}) \rangle = \left(\mathbf{x}^{(i)T} \mathbf{x}^{(j)} \right)^d$$

- Gaussian kernel / RBF kernel

$$\kappa(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \langle \phi(\mathbf{x}^{(i)}), \phi(\mathbf{x}^{(j)}) \rangle = \exp \left(-\frac{\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|^2}{2\sigma^2} \right)$$
$$\kappa(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \langle \phi(\mathbf{x}^{(i)}), \phi(\mathbf{x}^{(j)}) \rangle = \exp \left(-\gamma \|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|^2 \right)$$

Contoh

- Polinomial derajat 2 pada data latih 2 fitur / dimensi:

$$\phi(x_1, x_2) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$

- Dot product hasil feature mapping = kernel dari dot product vektor asal

$$\begin{aligned}\phi(\mathbf{a})^T \phi(\mathbf{b}) &= \begin{pmatrix} a_1^2 \\ \sqrt{2}a_1a_2 \\ a_2^2 \end{pmatrix}^T \begin{pmatrix} b_1^2 \\ \sqrt{2}b_1b_2 \\ b_2^2 \end{pmatrix} = a_1^2b_1^2 + 2a_1b_1a_2b_2 + a_2^2b_2^2 \\ &= (a_1b_1 + a_2b_2)^2 = \left(\begin{pmatrix} a_1 \\ a_2 \end{pmatrix}^T \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} \right)^2 = (\mathbf{a}^T \mathbf{b})^2\end{aligned}$$

$$\begin{pmatrix} a_1 \\ a_2 \end{pmatrix} \cdot \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}$$

Uji Pemahaman

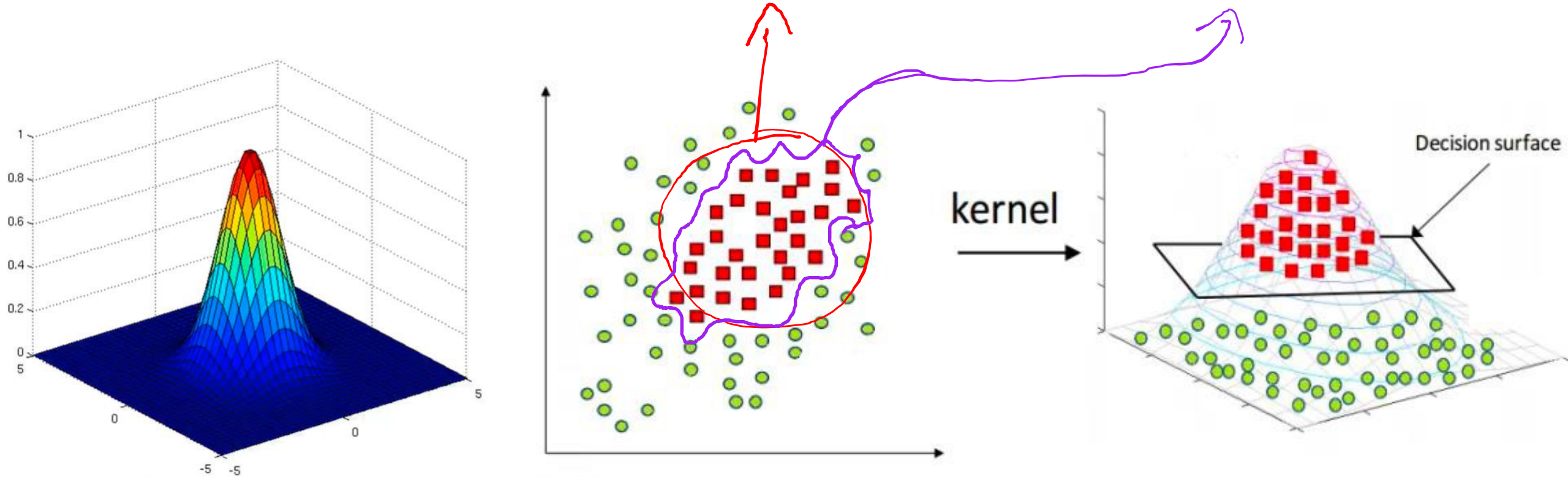
- $\mathbf{x}^{(1)} = (x_1^{(1)}, x_2^{(1)}, x_3^{(1)}) = (1, 2, 3)$, $\mathbf{x}^{(2)} = (x_1^{(2)}, x_2^{(2)}, x_3^{(2)}) = (4, 5, 6)$
- $\phi(\mathbf{x}^{(1)}) =$
 $(x_1^{(1)}x_1^{(1)}, x_1^{(1)}x_2^{(1)}, x_1^{(1)}x_3^{(1)}, x_2^{(1)}x_1^{(1)}, x_2^{(1)}x_2^{(1)}, x_2^{(1)}x_3^{(1)}, x_3^{(1)}x_1^{(1)}, x_3^{(1)}x_2^{(1)}, x_3^{(1)}x_3^{(1)}) =$
 $(1, 2, 3, 2, 4, 6, 3, 6, 9)$
- $\phi(\mathbf{x}^{(2)}) =$
 $(x_1^{(2)}x_1^{(2)}, x_1^{(2)}x_2^{(2)}, x_1^{(2)}x_3^{(2)}, x_2^{(2)}x_1^{(2)}, x_2^{(2)}x_2^{(2)}, x_2^{(2)}x_3^{(2)}, x_3^{(2)}x_1^{(2)}, x_3^{(2)}x_2^{(2)}, x_3^{(2)}x_3^{(2)}) =$
 $(16, 20, 24, 20, 25, 30, 24, 30, 36)$
- Hitung $\kappa(\mathbf{x}^{(1)}, \mathbf{x}^{(2)})$ secara manual, lalu bandingkan dengan kernel tricks

$$\phi(\mathbf{x}^{(1)}) \cdot \phi(\mathbf{x}^{(2)}) = 16 + 40 + 72 + 40 + 100 + 180 + 72 + 180 + 324 = 1024$$

$$\phi(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) = \phi(1, 2, 3, 4, 5, 6) = 16 + 100 + 324 + 2 \times 40 + 2 \times 180 + 2 \times 72 = 1024$$

RBF Kernel

- $\kappa(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \langle \phi(\mathbf{x}^{(i)}), \phi(\mathbf{x}^{(j)}) \rangle = \exp\left(-\gamma \|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|^2\right)$
- $\gamma \rightarrow$ hyperparameter: γ terlalu kecil \rightarrow underfitting, γ terlalu besar \rightarrow overfitting



Optimisasi dengan Kernel

- Proses optimisasi dapat dilihat sebagai primal maupun dual problem ([https://en.wikipedia.org/wiki/Duality \(optimization\)](https://en.wikipedia.org/wiki/Duality_(optimization)))
- Dalam optimisasi SVM secara dual, hipotesis akhir dapat ditulis sebagai:

$$\underline{w} = \sum_{i=1}^N \alpha_i y^{(i)} \Phi(\mathbf{x}^{(i)})$$
$$f(x) = \mathbf{w}^T \Phi(\mathbf{x}) + b = \sum_{i=1}^N \alpha_i y^{(i)} \Phi(\mathbf{x}^{(i)}) \Phi(\mathbf{x}) + b = \sum_{i=1}^N \alpha_i y^{(i)} \kappa(\mathbf{x}_i, \mathbf{x}) + b$$

Handwritten red annotations: An arrow points from $\Phi(\mathbf{x}^{(i)})$ in the first equation to $\Phi(\mathbf{x}^{(i)})$ in the second equation. A red box encloses the final term $\sum_{i=1}^N \alpha_i y^{(i)} \kappa(\mathbf{x}_i, \mathbf{x}) + b$, with an arrow pointing to the word "kernel" written above it.

Kernel Learning

- Manakah kernel terbaik dari banyak alternatif kernel?
- Multiple kernel learning (MKL): mencari kombinasi optimal dari beberapa kernel

$$\kappa(\mathbf{a}, \mathbf{b}, ; \gamma) = \sum_{k=1}^m \gamma_k \kappa_k(\mathbf{a}, \mathbf{b})$$

- MKL framework → konstruksi fungsi kernel:

- Berbagai macam fitur

- Berbagai tipe kernel

- Berbagai parameter kernel

→ RBF, poly 2, poly 3, ...
 γ, d

Curse of kernelization

- Melatih kernel classifier sangat computationally expensive
- SVM linear butuh linear time $O(N)$
- SVM kernel dengan QP solvers butuh $O(N^3)$
- Bagaimana cara melatih pada large-scale dataset? Gunakan aproksimasi kernel

Aproksimasi Kernel

- Tujuan: buat representasi baru $\mathbf{z}(\mathbf{x}) \in \mathbb{R}^D$ sehingga $\kappa(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) \approx \mathbf{z}(\mathbf{x}^{(i)})^T \mathbf{z}(\mathbf{x}^{(j)})$
- Model linear:
 - Hipotesis dapat ditulis ulang sebagai:

$$f(x) = \sum_{i=1}^M \alpha_i \kappa(\mathbf{x}^{(i)}, \mathbf{x}) \approx \sum_{i=1}^M \alpha_i \mathbf{z}(\mathbf{x}^{(i)})^T \mathbf{z}(\mathbf{x}) = \mathbf{w}^T \mathbf{z}(\mathbf{x})$$

dimana $\mathbf{w}^T = \sum_{i=1}^M \alpha_i \mathbf{z}(\mathbf{x}^{(i)})$

- Gunakan linear classifier pada representasi baru \mathbf{z}
- 2 metode:
 - Aproksimasi fungsional: metode Fourier
 - Aproksimasi matriks: metode Nystrom

SVM Multi Kelas



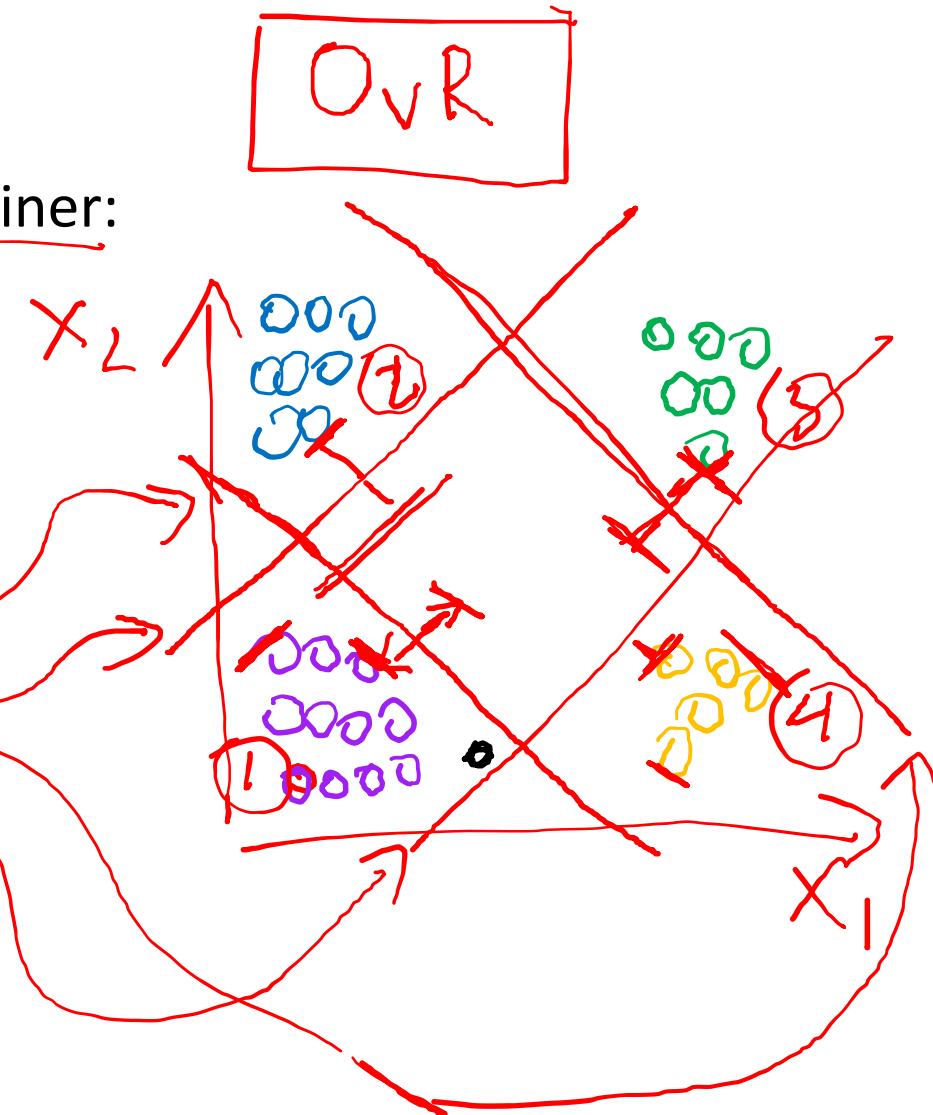
God's People for God's Glory

CALVIN
INSTITUTE OF TECHNOLOGY

Klasifikasi multi kelas

- Misalkan terdapat K kelas
- Gunakan one versus rest (OvR): latih sebanyak K SVM biner:
 - 1st class vs (2-k)th class
 - 2nd class vs (1,3-k)th class
 - ...
- K decision functions:

$$\begin{aligned} &(\mathbf{w}_1)^T \phi(\mathbf{x}) + b_1 \\ &(\mathbf{w}_2)^T \phi(\mathbf{x}) + b_2 \\ &\vdots \\ &(\mathbf{w}_K)^T \phi(\mathbf{x}) + b_K \end{aligned}$$



OvR

- Prediksi

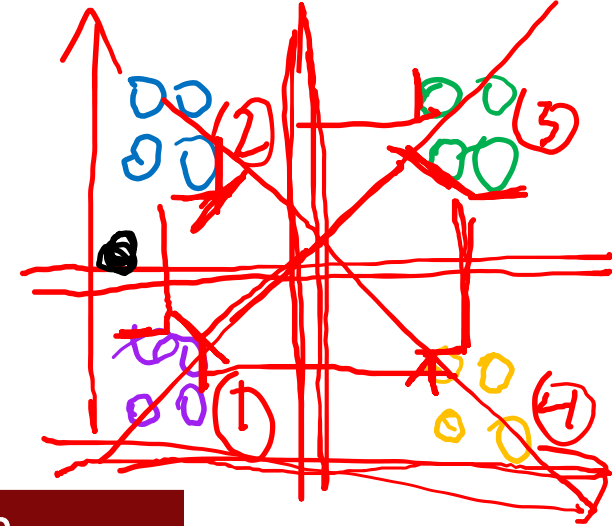
$$\arg \max_k (\mathbf{w}_k)^T \phi(\mathbf{x}) + b_k$$

- Jika prediksinya adalah 1st class:

$$\begin{aligned} (\mathbf{w}_1)^T \phi(\mathbf{x}) + b_1 &\geq +1 \\ (\mathbf{w}_2)^T \phi(\mathbf{x}) + b_2 &\leq -1 \\ &\vdots \\ (\mathbf{w}_K)^T \phi(\mathbf{x}) + b_K &\leq -1 \end{aligned}$$

$$\frac{4(4-1)}{2} = 6$$

- One versus one: latih $\frac{K(K-1)}{2}$ SVM biner:
 $(1,2), (1,3), \dots, (1,k), (2,3), (2,4), \dots, (k-1,k)$
- Contoh: jika terdapat 4 kelas \rightarrow butuh 6 SVM biner



$y^{(i)} = +1$	$y^{(i)} = -1$	Fungsi keputusan
Kelas 1	Kelas 2	$f_{12}(\mathbf{x}) = (\mathbf{w}_{12})^T \mathbf{x} + b_{12}$
Kelas 1	Kelas 3	$f_{13}(\mathbf{x}) = (\mathbf{w}_{13})^T \mathbf{x} + b_{13}$
Kelas 1	Kelas 4	$f_{14}(\mathbf{x}) = (\mathbf{w}_{14})^T \mathbf{x} + b_{14}$
Kelas 2	Kelas 3	$f_{23}(\mathbf{x}) = (\mathbf{w}_{23})^T \mathbf{x} + b_{23}$
Kelas 2	Kelas 4	$f_{24}(\mathbf{x}) = (\mathbf{w}_{24})^T \mathbf{x} + b_{24}$
Kelas 3	Kelas 4	$f_{34}(\mathbf{x}) = (\mathbf{w}_{34})^T \mathbf{x} + b_{34}$

- Pilih kelas dengan voting terbanyak

OvR vs OvO

- Asumsikan optimisasi SVM dengan ukuran n adalah $O(n^d)$
- OvR: terdapat K model, masing-masing n data $\rightarrow K O(n^d)$
- OvO: terdapat $\frac{K(K-1)}{2}$ model, masing-masing $2n/K$ data $\rightarrow \frac{K(K-1)}{2} O((2n/K)^d)$

$\rightarrow QP \rightarrow d=3$

$$OvR : 10 \times O(1M^3)$$

$$OvO : 45 \times O(200K^3)$$

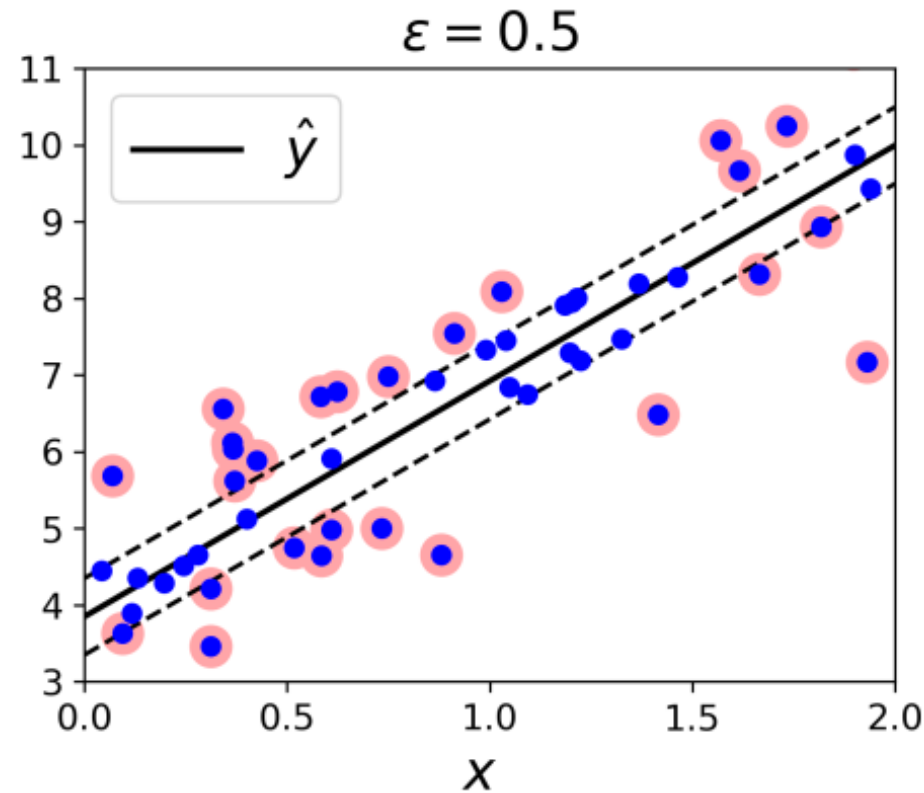
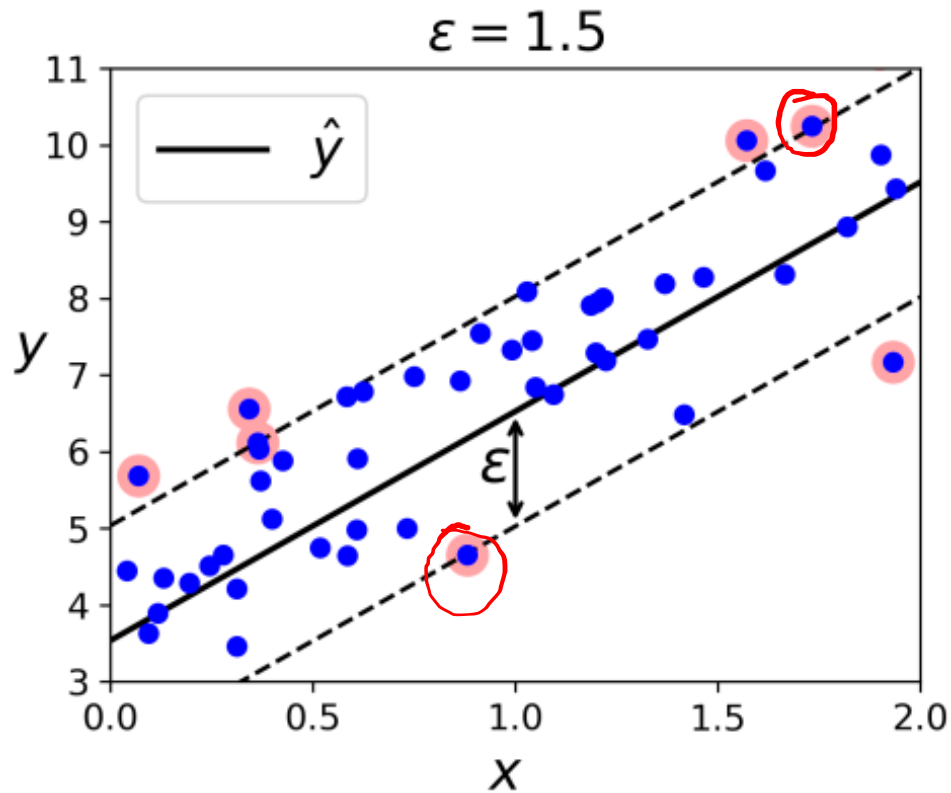
Regresi SVM

Regresi dengan SVM

- SVM dapat digunakan untuk regresi
- Klasifikasi:
 - membuat margin selebar mungkin di antara dua kelas
 - Meminimalkan margin violation: sesedikit mungkin data yang melewati margin
- Regresi:
 - Membuat sebanyak mungkin data masuk dalam margin
 - Meminimalkan margin violation: sesedikit mungkin data melewati margin

Regresi SVM

- Lebar margin dikontrol oleh hyperparameter ϵ
- Apa akibat penambahan data di dalam margin?
- Tidak berdampak pada prediksi model $\rightarrow \epsilon$ -insensitive (linear)



Optimisasi Regresi SVM

- Tujuan: memperoleh model $\hat{f}(x)$ dengan prediksi \hat{y} yang berbeda dengan target y maksimal sejauh ϵ
- Formulasi permasalahan optimisasi dengan soft constraint:

$$\min_{w,b,\zeta,\zeta^*} \left[\frac{1}{2} w^T w \right] + C \left[\sum_{i=1}^n (\zeta_i + \zeta_i^*) \right]$$

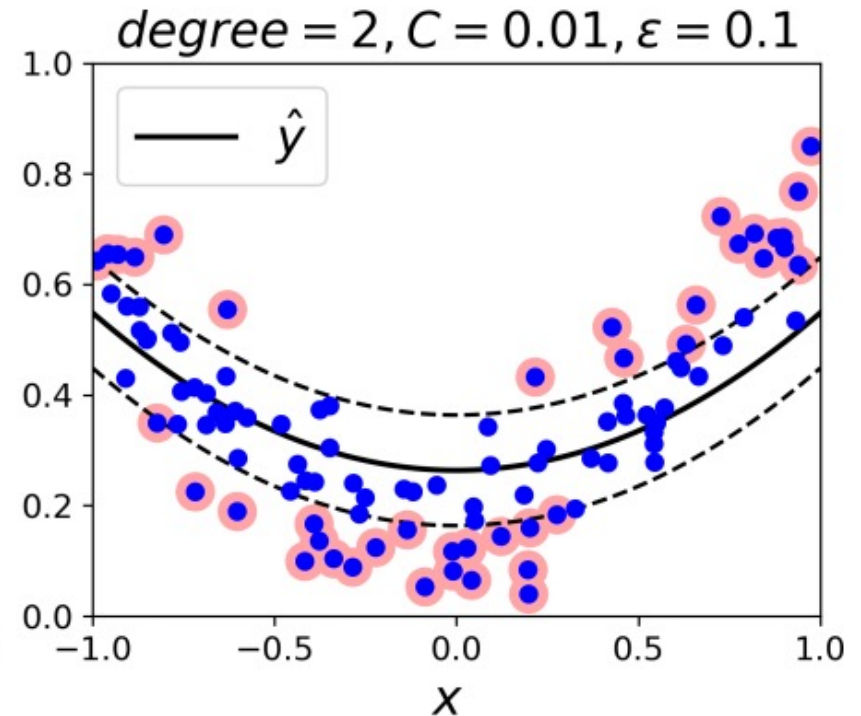
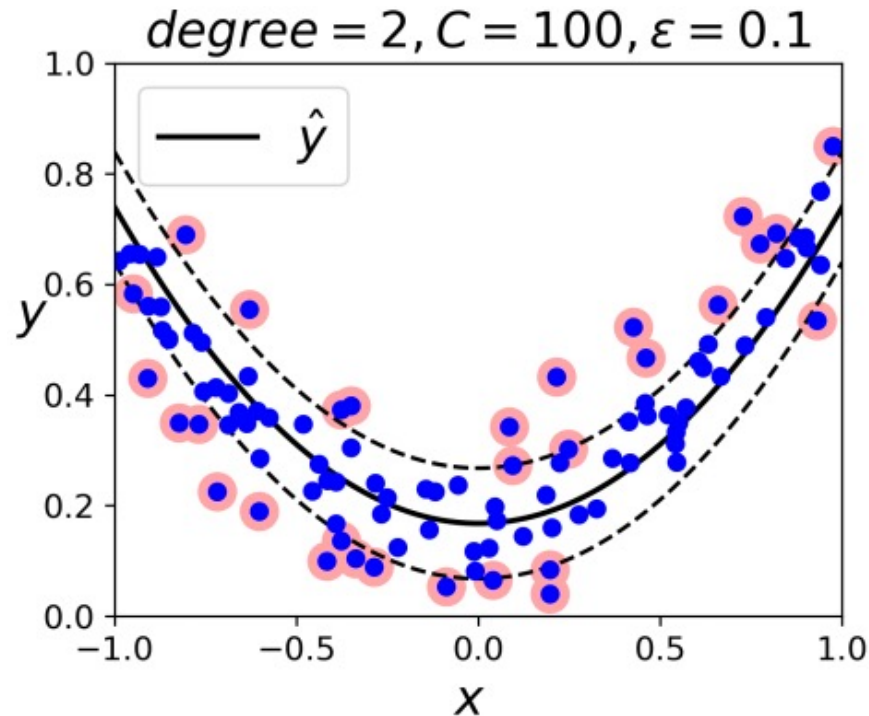
margin (pointing to $\frac{1}{2} w^T w$) *Violation* (pointing to $\sum_{i=1}^n (\zeta_i + \zeta_i^*)$)

$$\begin{aligned} \text{subject to } & y_i - w^T \phi(x_i) - b \leq \epsilon + \zeta_i, \\ & w^T \phi(x_i) + b - y_i \leq \epsilon + \zeta_i^*, \\ & \zeta_i, \zeta_i^* \geq 0, i = 1, \dots, n \end{aligned}$$

Regresi SVM Nonlinear

- Untuk data non-linear, kita dapat memakai SVM dengan kernel
- Ilustrasi: data kuadratik (polynomial derajat 2) dengan hyperparameter C yang berbeda
- Kasus mana lebih underfitting?

Underfitting



Tuhan Memberkati



God's People for God's Glory

CALVIN
INSTITUTE OF TECHNOLOGY