

# Support Vector Machines

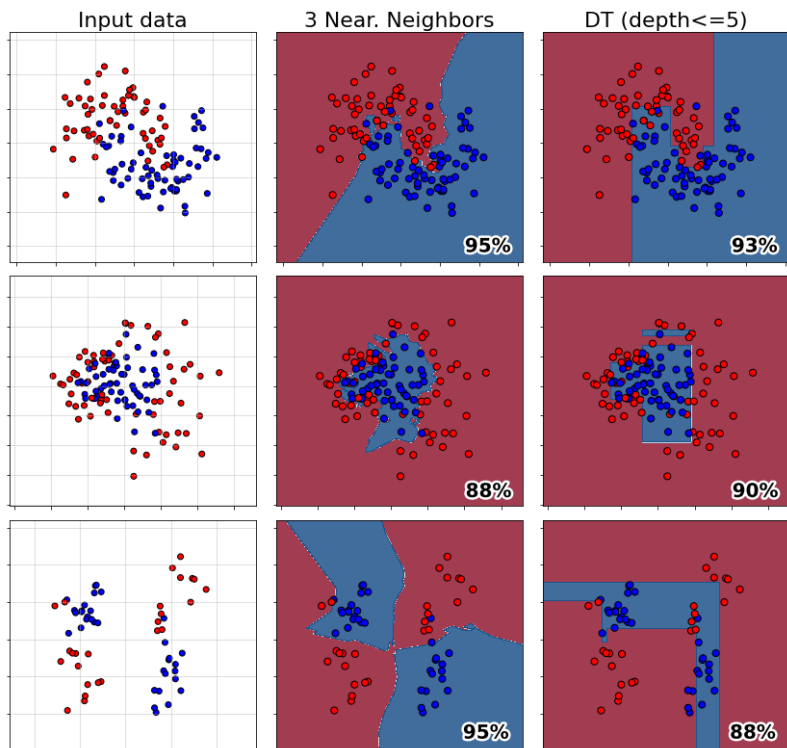
---

# Outline

- Hard SVM
- Soft SVM
- Kernel SVM

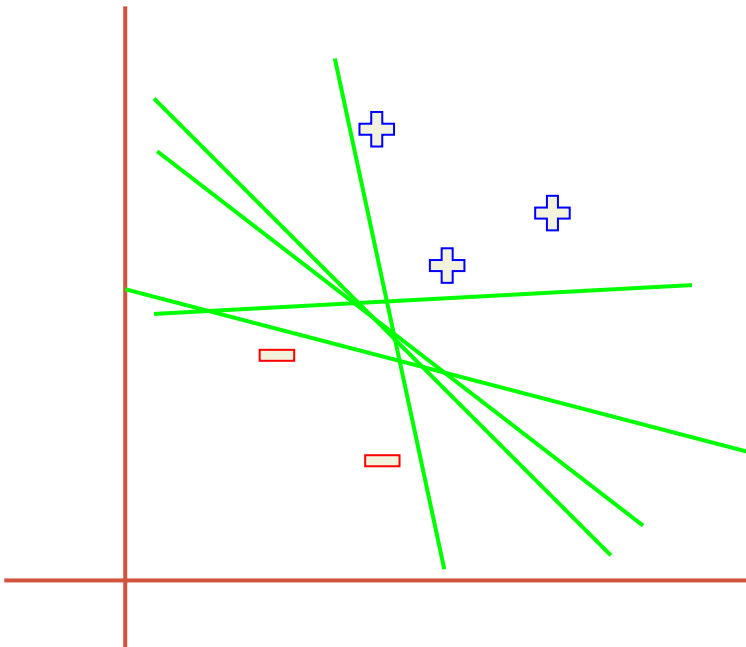
# Decision boundaries

- Different models on 3 datasets.
- On the bottom right – the train accuracy.

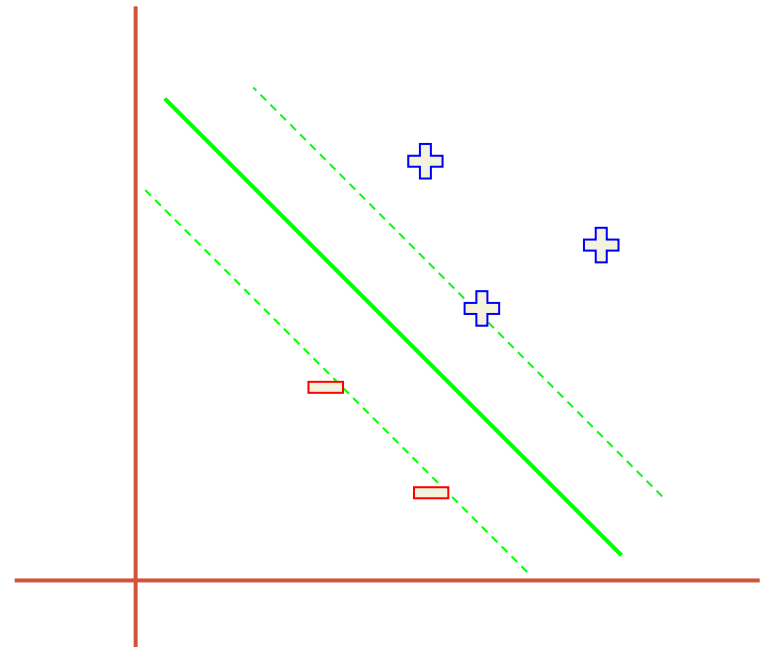


# Linear separation

Infinitely many separators



Max-margin separator



Intuition: Hard SVM fits the “widest” possible strip between classes

# Hard SVM

---

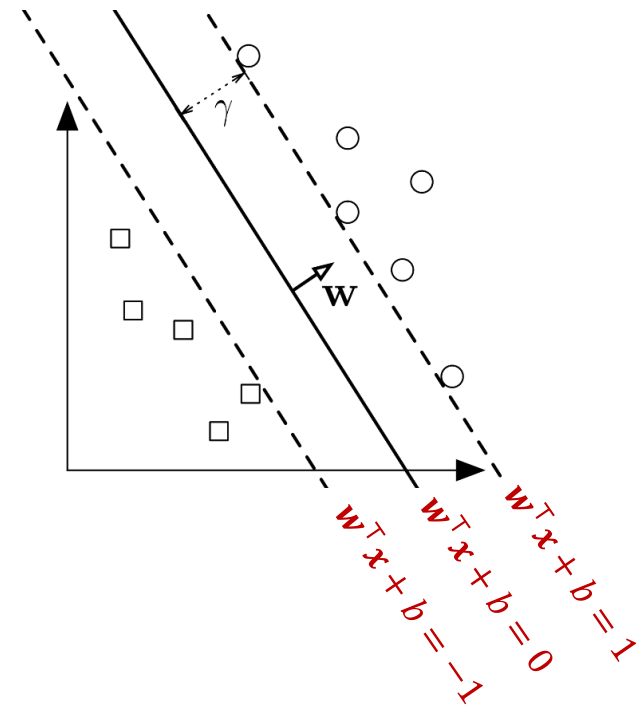
# Hard SVM: Recap

- We wanted a vector  $\mathbf{w}$  that maximizes the margin.

$$\begin{aligned} & \text{margin, } \gamma(\mathbf{w}; S) \\ \operatorname{argmax}_{\mathbf{w} \in \mathbb{R}^d} & \min_{i \in [m]} \frac{|\mathbf{w}^\top \mathbf{x}_i|}{\|\mathbf{w}\|_2} \\ \text{s.t. } & y_i \cdot \mathbf{w}^\top \mathbf{x}_i \geq 0, \quad \forall i \in [m] \end{aligned}$$

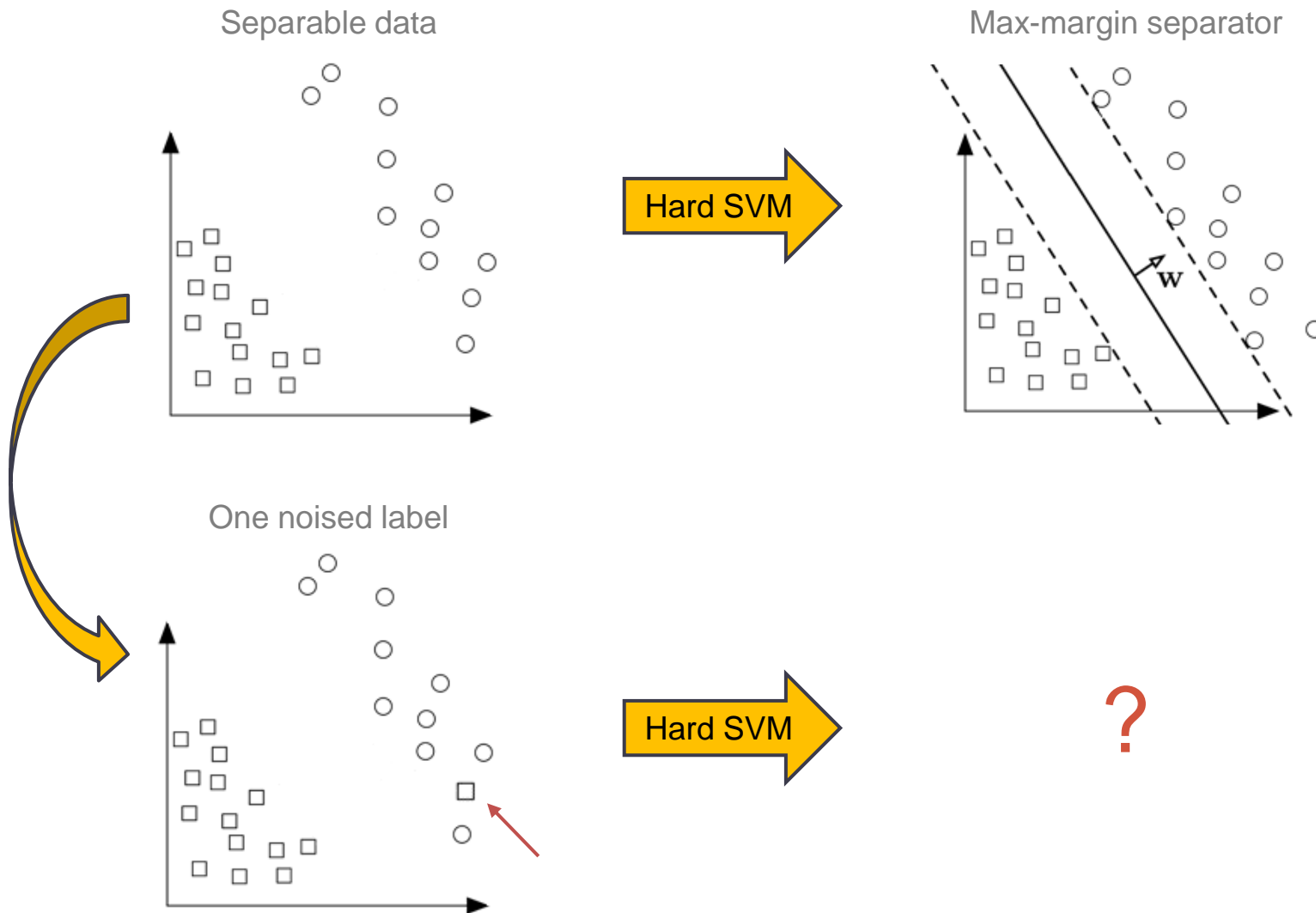
- We solved an equivalent optimization problem:

$$\begin{aligned} \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} & \|\mathbf{w}\|_2^2 \\ \text{s.t. } & y_i \cdot \mathbf{w}^\top \mathbf{x}_i \geq 1, \quad \forall i \in [m] \end{aligned}$$



Classify all points correctly  
with a margin at least  $1/\|\mathbf{w}\|_2^2$

# Sensitivity to outliers



# Issues with Hard SVM

- Only works with linearly separable data
- Highly sensitive to outliers
- More flexible – **Soft SVM**:
  - Balances between **max-margin** and **margin violations**.



# Soft SVM

---

# Soft SVM: Optimization problem

- Two conflicting objectives:

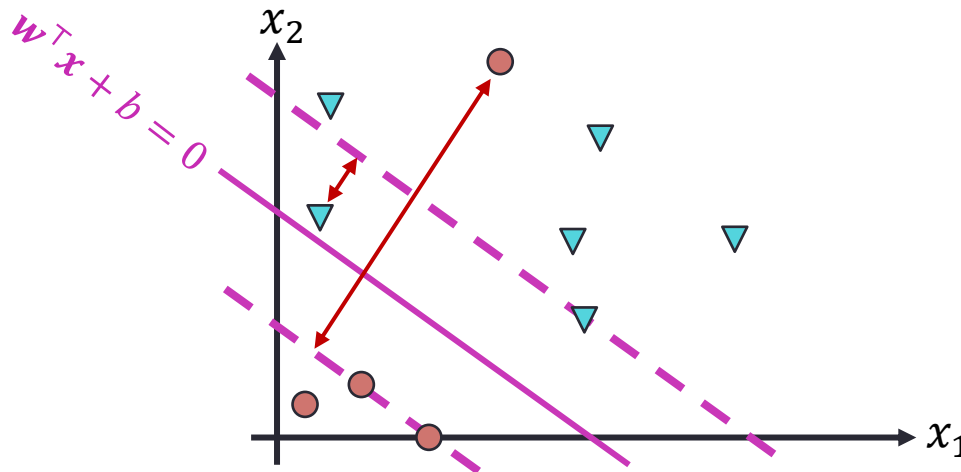
$$\operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} \underbrace{\lambda \|\mathbf{w}\|_2^2}_{\text{Complexity penalty (Regularization)}} + \frac{1}{m} \sum_{i \in [m]} \underbrace{\max\{0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b)\}}_{\text{Margin violation penalty}}$$

Complexity penalty (Regularization)

Hyperparameter  $\lambda$  controls the model complexity

Margin violation penalty

How much the  $i^{\text{th}}$  example violates the margin



# Soft SVM: Optimization problem

- Two conflicting objectives:

$$\operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} \lambda \|\mathbf{w}\|_2^2 + \frac{1}{m} \sum_{i \in [m]} \max\{0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b)\}$$

- Equivalently (convince yourselves), some sources formulate as:

$$\operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} \|\mathbf{w}\|_2^2 + c \sum_{i \in [m]} \max\{0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b)\}$$

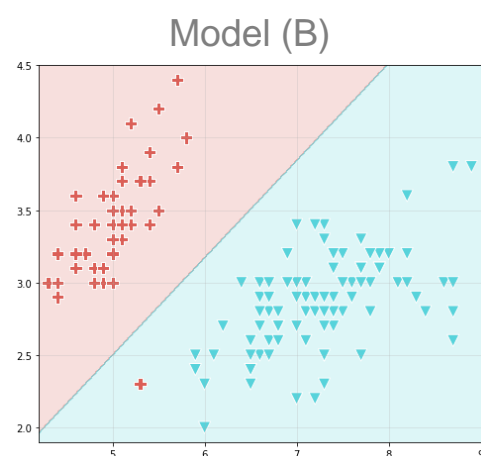
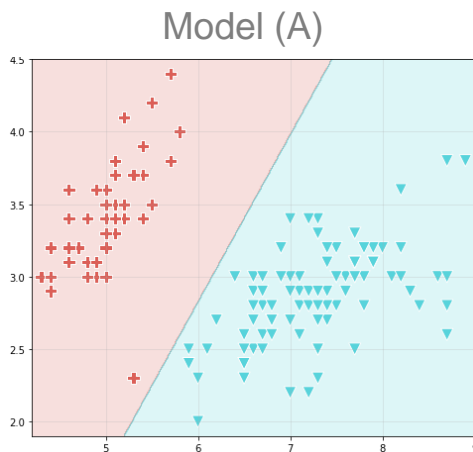
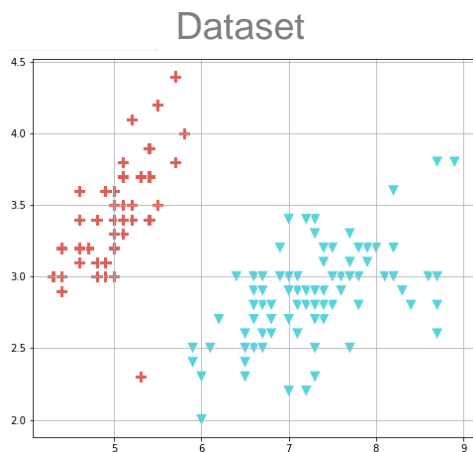
# Soft SVM: Separable case

$$\operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} \|\mathbf{w}\|_2^2 + C \sum_{i \in [m]} \max\{0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b)\}$$

- Given a linearly separable dataset, we train two models with:  $C=1$  and  $1000$ .

```
from sklearn.svm import LinearSVC
models = [LinearSVC(C=1, max_iter=10000), LinearSVC(C=1000, max_iter=10000)]
models = [h.fit(X, y) for h in models]
```

- Exercise:** match the models to  $C=1$  and  $1000$ .



# Tuning the regularizer $\lambda$

- $\lambda$  balances between model complexity and the violation penalty.

$$\operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} \underbrace{\lambda \|\mathbf{w}\|_2^2}_{\text{Complexity penalty (Regularization)}} + \frac{1}{m} \sum_{i \in [m]} \underbrace{\max\{0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b)\}}_{\text{Margin violation penalty}}$$

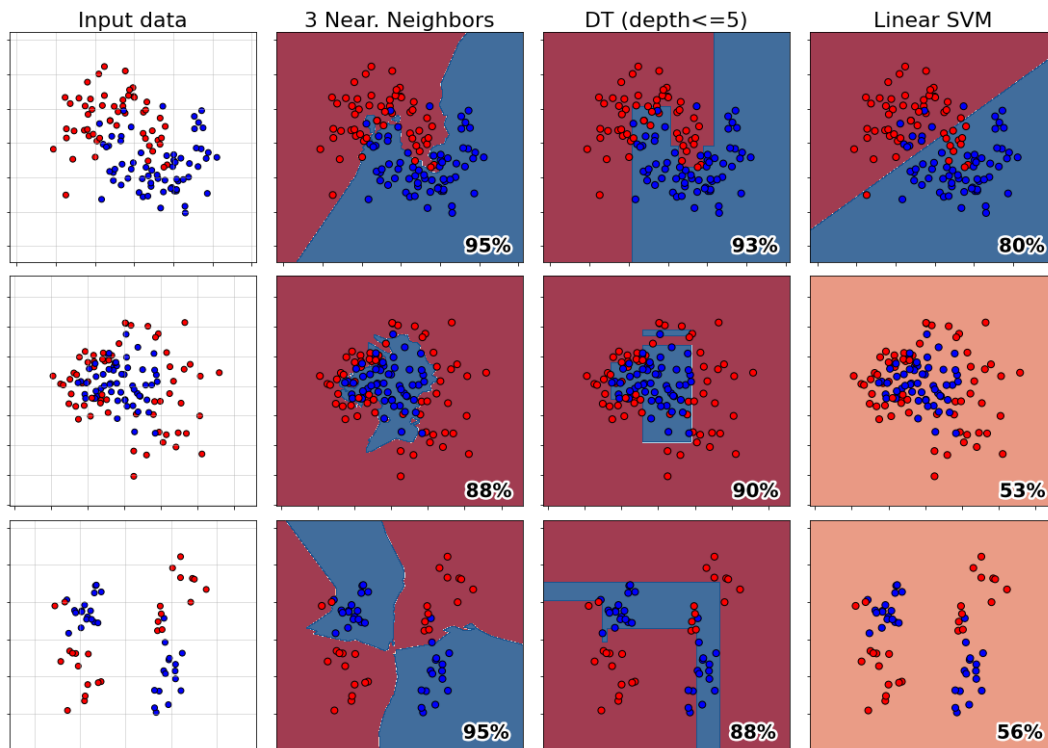
Hyperparameter  $\lambda$  controls the model complexity

How much the  $i^{\text{th}}$  example *violates* the margin

- Larger  $\lambda$ : **more tolerance** to violations  $\Rightarrow$  **lower complexity**
- Smaller  $\lambda$ : **less tolerance** to violations  $\Rightarrow$  **overfitting**
- As  $\lambda$  decreases, we get closer to Hard-SVM solution. **Why?**  
(if data is separable)

# Decision boundaries

- Different models on 3 datasets.
- On the bottom right – the train accuracy.

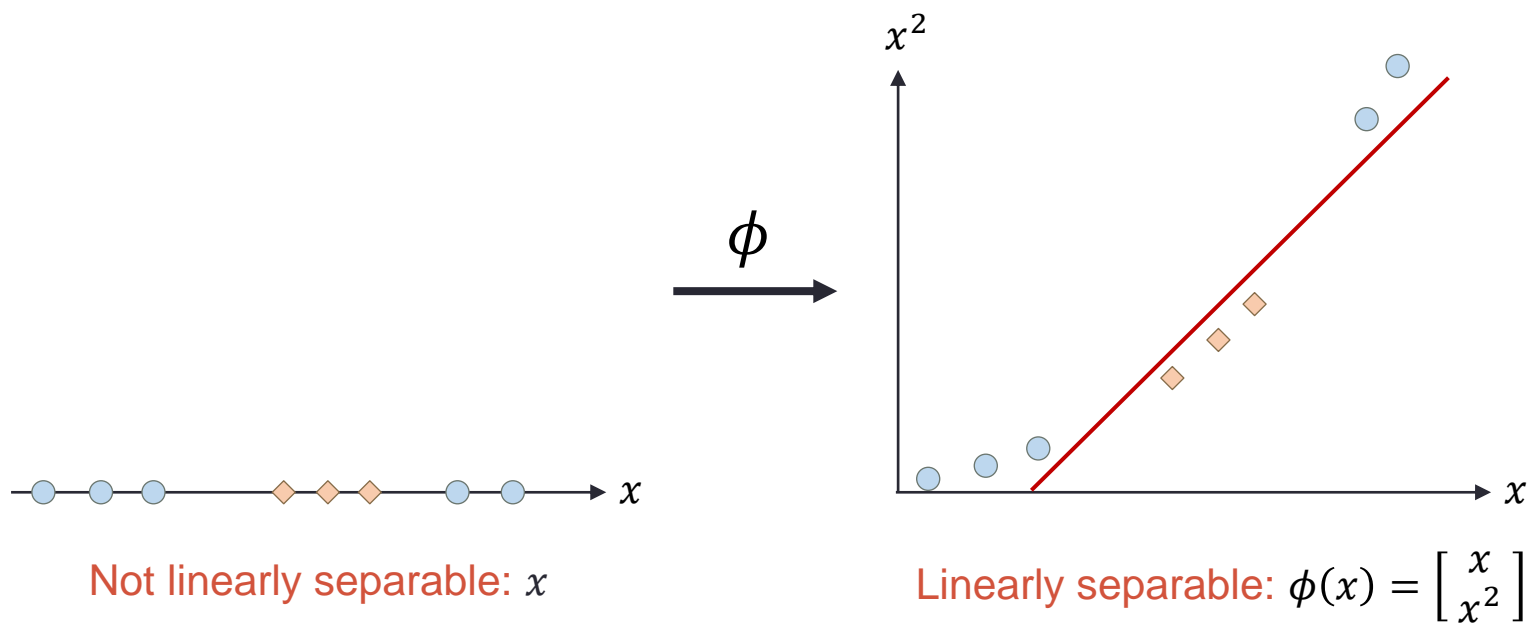


# Kernel SVM

---

# Feature mappings

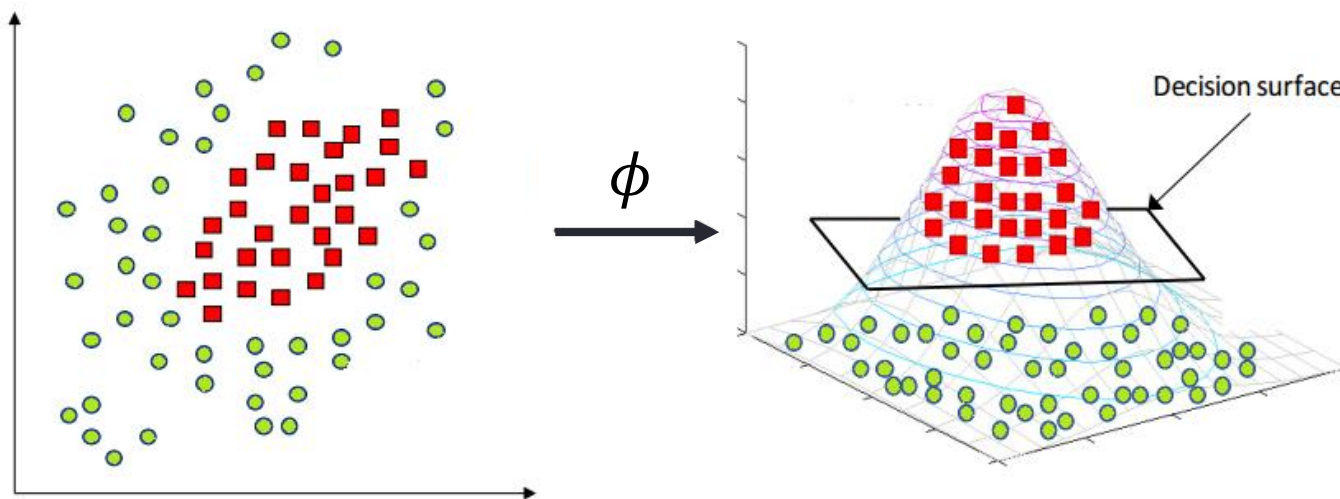
- Linear SVMs are efficient and often work well
- However, many problems are **not linearly separable**
  - One approach is to use **a feature mapping  $\phi$**  to add more features:





# Feature mappings

- Linear SVMs are efficient and often work well
- However, many problems are **not linearly separable**
  - One approach is to use **a feature mapping  $\phi$**  to add more features:



# The polynomial mapping

- The 2<sup>nd</sup>-deg. polynomial mapping:  $\phi \left( \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \right) = \begin{bmatrix} 1 \\ \sqrt{2}u_1 \\ \sqrt{2}u_2 \\ \sqrt{2}u_1u_2 \\ u_1^2 \\ u_2^2 \end{bmatrix}$

- Creates decision boundaries of the form:

$$0 \leq \mathbf{w}^\top \phi \left( \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right) = w_1 + \sqrt{2}w_2x_1 + \sqrt{2}w_3x_2 + \sqrt{2}w_4x_1x_2 + w_5x_1^2 + w_6x_2^2$$

- Let us visualize it

What shapes can we create?  
linear, circles, ellipsis, hyperbolas

# The polynomial mapping

- The 2<sup>nd</sup>-deg. **polynomial mapping**:  $\phi \left( \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \right) = \begin{bmatrix} 1 \\ \sqrt{2}u_1 \\ \sqrt{2}u_2 \\ \sqrt{2}u_1u_2 \\ u_1^2 \\ u_2^2 \end{bmatrix}$

- Creates **decision boundaries** of the form:

$$0 \leq \mathbf{w}^\top \phi \left( \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right) = w_1 + \sqrt{2}w_2x_1 + \sqrt{2}w_3x_2 + \sqrt{2}w_4x_1x_2 + w_5x_1^2 + w_6x_2^2$$

- Let us visualize it
- Works great with ML algorithms
- Higher polynomials can fit more complex data
  - But more features might cause **computational** (and statistical) problems

# The kernel trick

## Hard-SVM

$$\begin{aligned} \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{w}\|_2^2 \\ \text{s.t. } y_i \cdot \mathbf{w}^\top \mathbf{x}_i \geq 1, \quad \forall i \in [m] \end{aligned}$$

Feature mapping



$$\begin{aligned} \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^{d'}} \|\mathbf{w}\|_2^2 \\ \text{s.t. } y_i \cdot \mathbf{w}^\top \phi(\mathbf{x}_i) \geq 1, \quad \forall i \in [m] \end{aligned}$$

- If  $d' \gg d$ , optimization is expensive.

## Dual problem uses only inner products

$$\max_{\alpha \in \mathbb{R}_+^m} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y_i y_j \alpha_i \alpha_j \mathbf{x}_i^\top \mathbf{x}_j$$

Kernel trick



$$\max_{\alpha \in \mathbb{R}_+^m} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y_i y_j \alpha_i \alpha_j \underbrace{\phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)}_{=K(\mathbf{x}_i, \mathbf{x}_j)}$$

- **The kernel trick:** solve SVM for any feature mapping  $\phi$ , by simply computing  $K(\mathbf{x}_i, \mathbf{x}_j)$  for every two training samples.

# Valid kernels

- A kernel must hold  $K(\mathbf{u}, \mathbf{v}) = \phi(\mathbf{u})^\top \phi(\mathbf{v})$  for some feature mapping  $\phi$ .

$$\mathbf{u}: \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} \quad \mathbf{v}: \begin{pmatrix} v_1 \\ v_2 \end{pmatrix} \quad K(\mathbf{u}, \mathbf{v}) = u_1^2 + 2u_1u_2 + u_2^2$$

- Prove:** the 2<sup>nd</sup>-deg. polynomial mapping  $\phi \left( \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \right) = \begin{bmatrix} 1 \\ \sqrt{2}u_1 \\ \sqrt{2}u_2 \\ \sqrt{2}u_1u_2 \\ u_1^2 \\ u_2^2 \end{bmatrix}$ 

$$\begin{matrix} 1 \\ 2u_1v_1 \\ 2u_2v_2 \\ 2u_1u_2v_1v_2 \\ u_1^2v_1^2 \\ u_2^2v_2^2 \end{matrix}$$

$$u^2v^2 + 2uv + 1$$

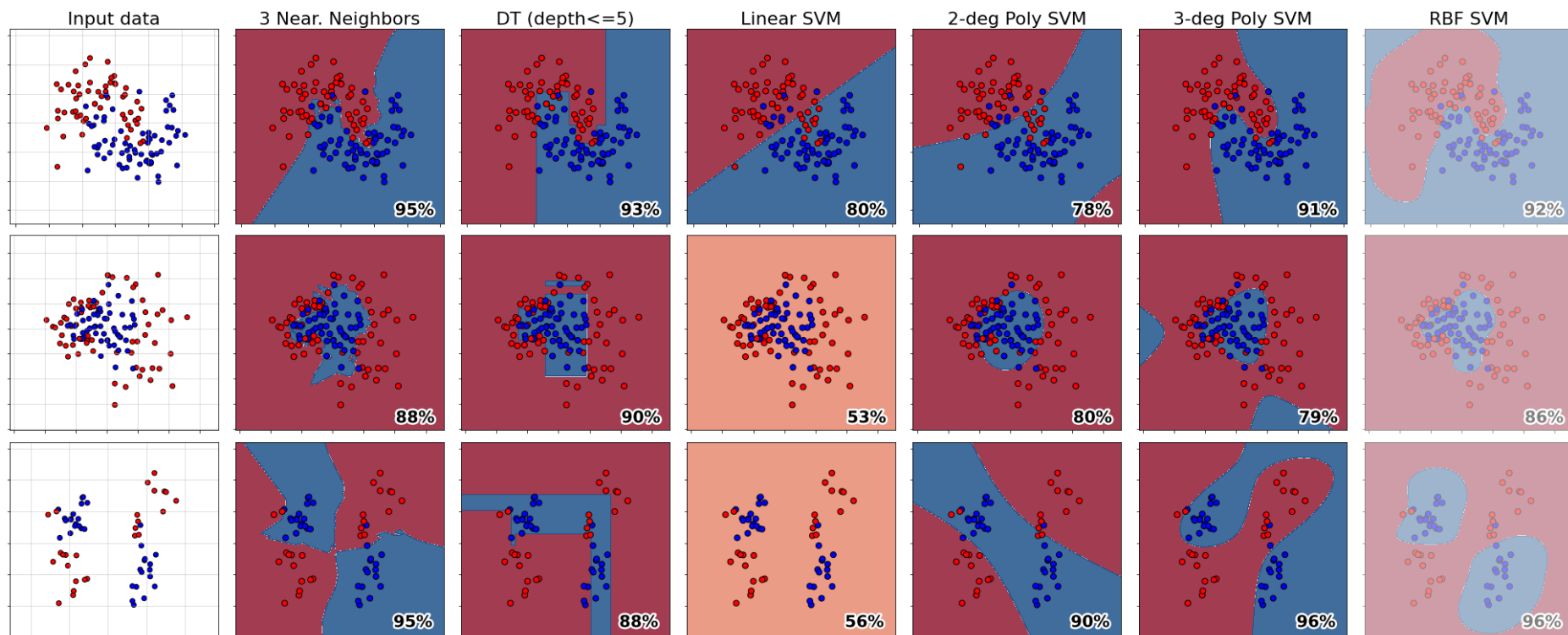
$$(u^2, \sqrt{2}u, 1)$$

defines a (polynomial) **kernel**  $K(\mathbf{u}, \mathbf{v}) = (\mathbf{u}^\top \mathbf{v} + 1)^2 = \phi(\mathbf{u})^\top \phi(\mathbf{v})$ .

- Computation shortcut:
  - Computing  $\phi(\mathbf{u}), \phi(\mathbf{v})$ :  $\mathcal{O}(d^2)$
  - Computing  $\phi(\mathbf{u})^\top \phi(\mathbf{v})$  directly:  $\mathcal{O}(d^2)$
  - Computing  $(\mathbf{u}^\top \mathbf{v} + 1)^2$ :  $\mathcal{O}(d)$
- Can now classify using all its monomials

# Decision boundaries

- Different models on 3 datasets.
- On the bottom right – the train accuracy.

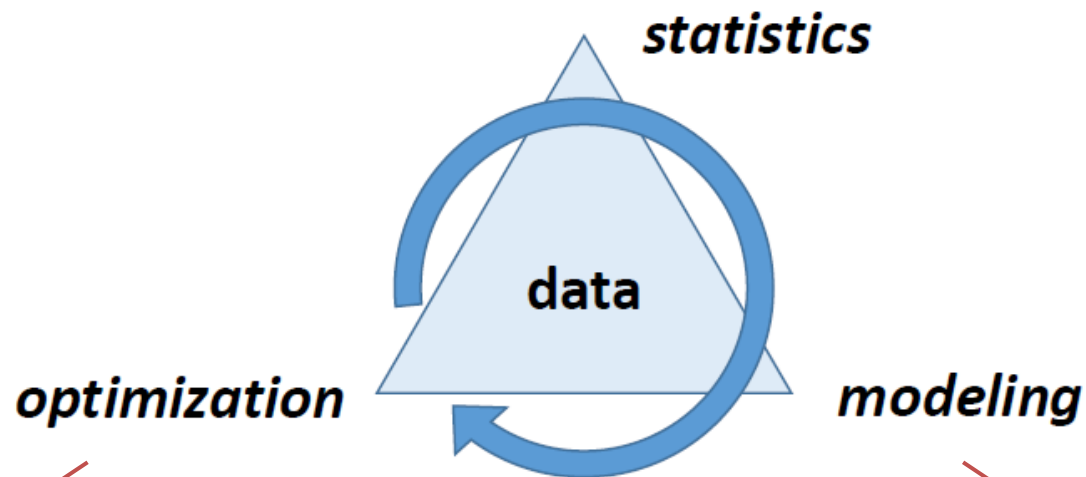


(in Major HW2)

# Which should we use?

- Rule of thumb: try the **linear kernel** first
  - Especially if training set is large (dual becomes expensive!)
  - In sklearn, `LinearSVC` is much faster than `SVC(kernel="linear")`
- If training set not too large, try **Gaussian RBF**
- If not satisfied, try others

# Summary



- Heavy optimization and math
  - Reread lecture if needed
- Global convergence
- Primal vs. Dual
- Struggles with higher dimensions and/or many samples

- SVM is flexible, high-performing
  - Linear and nonlinear
- Controlling model complexity
  - Through hyperparameter  $\lambda$  or  $C$
  - Kernel choice