# Support Vector Machines

**Justin Pounders**

# Introduction to SVM

Support Vector Machines (SVMs) are <span style="color:red">statistical models</span> used for **classification**.

Review:
- What is classification?
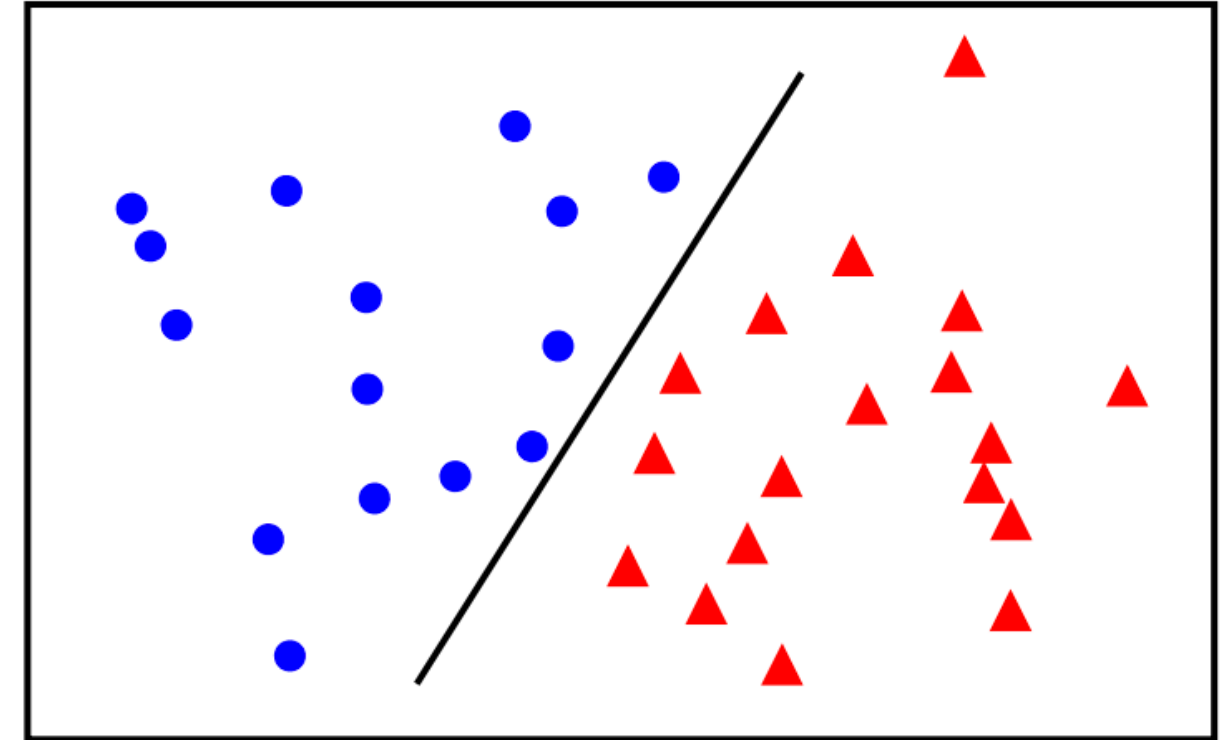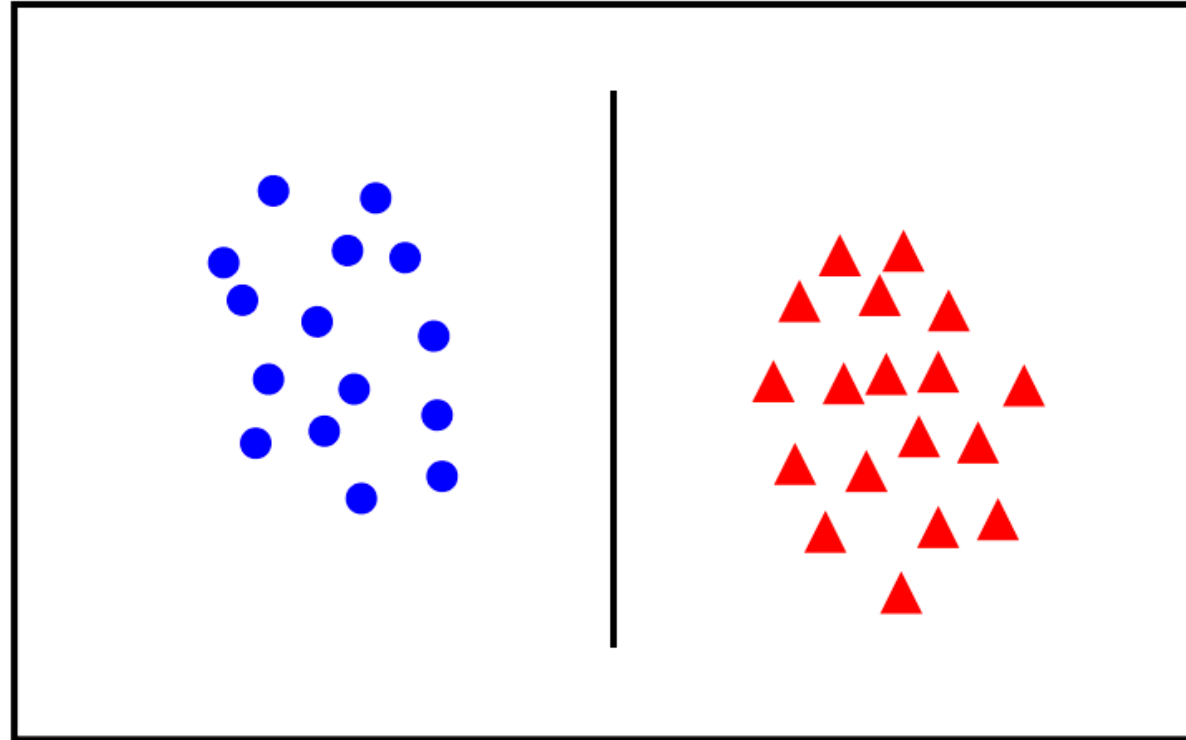- What other classification models have you seen?

# How to Think about SVM

— The geometric intuition of SVMs is easier to grasp than the

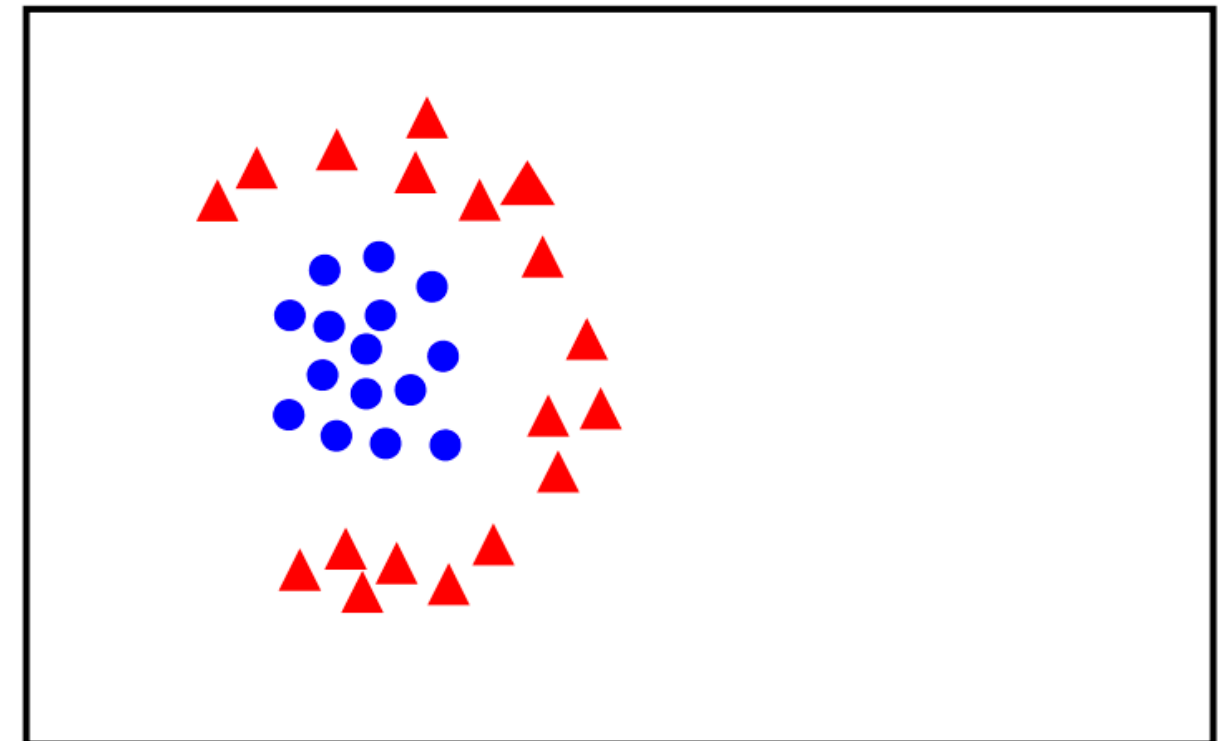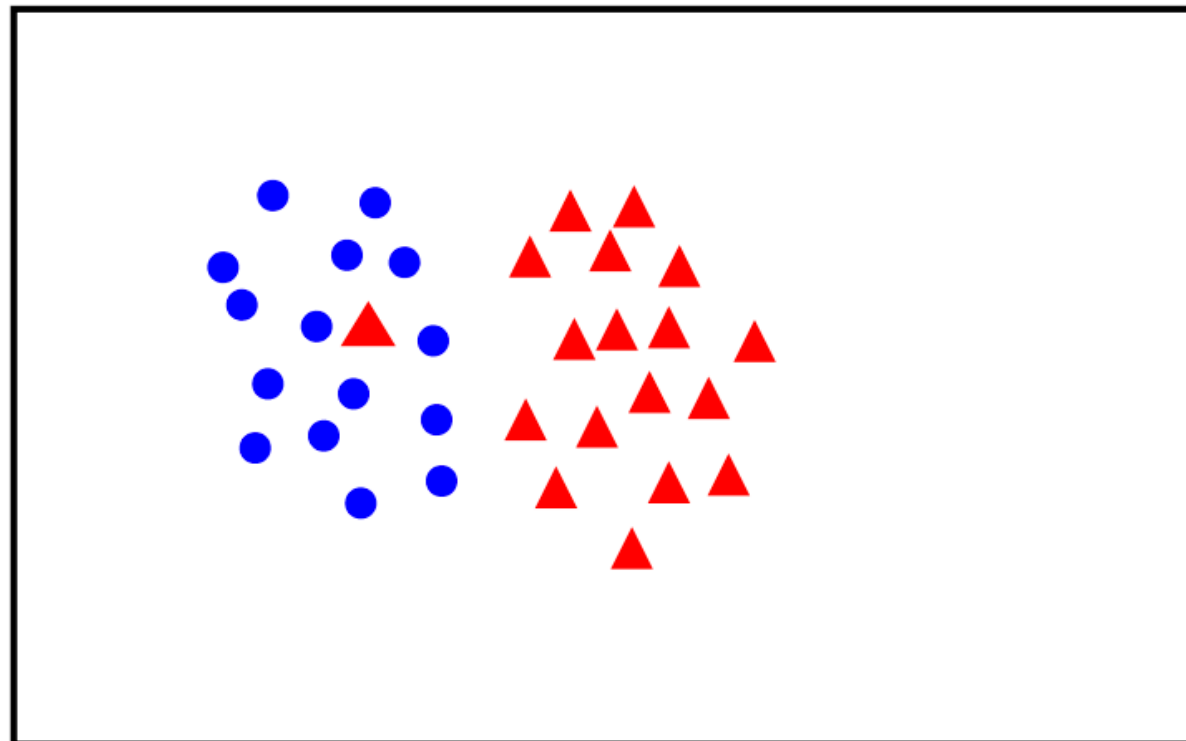— mathematical constructs needed to make it work

# Linear Separability

SVMs work really well for data in which the classes are linearly separable.
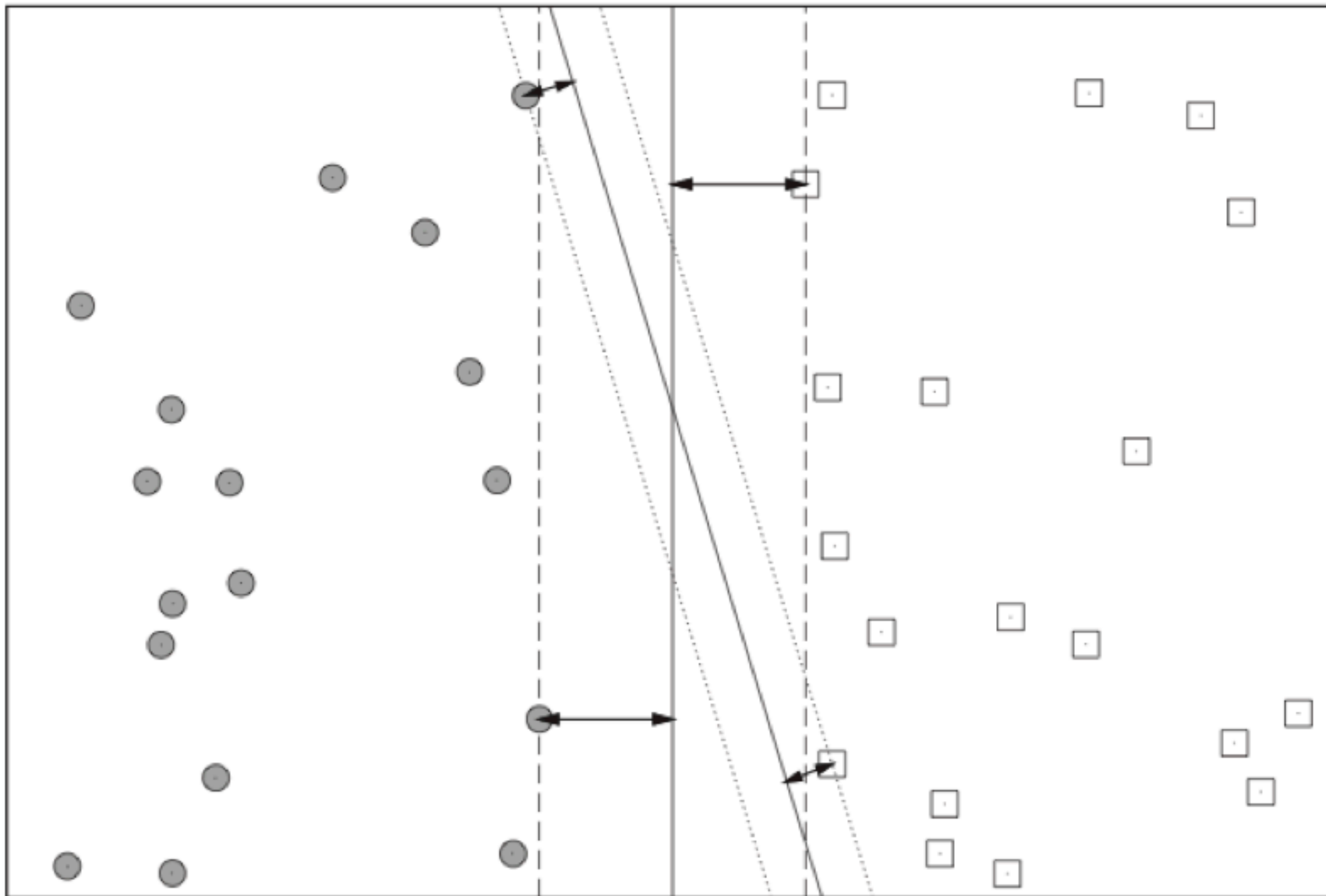
linearly
separable

not
linearly
separable

# Maximum-Margin Estimator

If classes are linearly seperable, SVM finds the hyperplane that seepearates the classess with maximum margin.
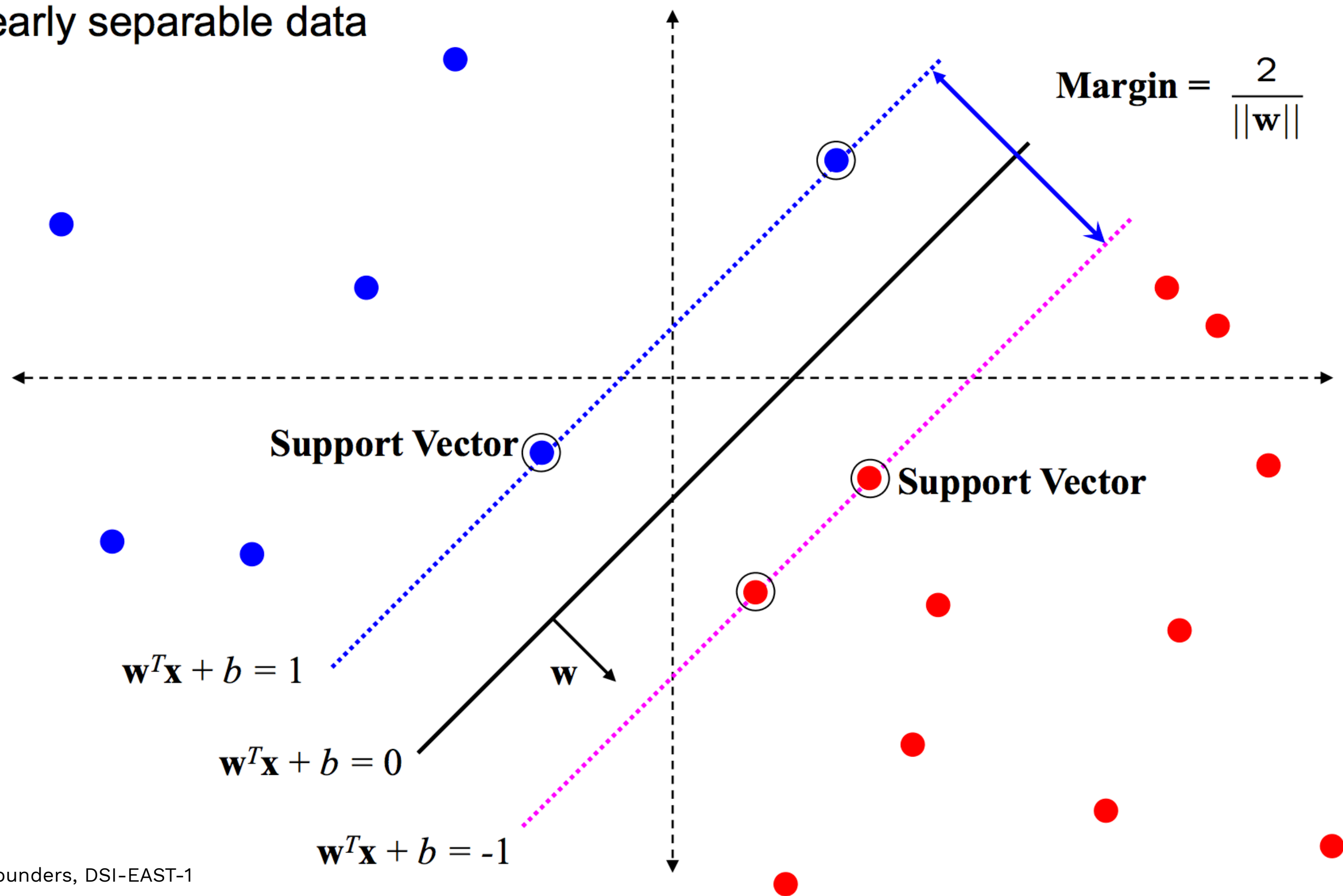
— **What is a hyperplane?**

**F I G U R E 18-4.** Two decision boundaries and their margins. Note that the vertical decision boundary has a wider margin than the other one. The arrows indicate the distance between the respective support vectors and the decision boundary.

# Why maximize the margin?

— SVM solves for a decision boundary that should minimize the generalization error.

— Observations near the decision boundary are the most "ambiguous"

— SVM defines it's fit using the most ambiguous points

# linearly separable data



Margin = $\dfrac{2}{||\mathbf{w}||}$

**Support Vector**

**Support Vector**

$\mathbf{w}^T\mathbf{x} + b = 1$

$\mathbf{w}$

$\mathbf{w}^T\mathbf{x} + b = 0$

$\mathbf{w}^T\mathbf{x} + b = -1$

# Maximum-Margin Hyperplane

**Goal**: Find $w$ that leads to the max-margin hyperplane

$$w \leftarrow \max_{w} \frac{2}{\|w\|} = \max \text{ margin}$$

subject to all points being on the "right side"

$$w^T x_i + b \geq 1 \quad \text{if } y_i = 1$$
$$w^T x_i + b \leq -1 \quad \text{if } y_i = -1$$

# Maximum-Margin Hyperplane

What if data are not linearly seperable?

# Maximum-Margin Hyperplane

What if data are not linearly seperable?

— Still want to minimize $\|w\|$ (maximize margin)

# Maximum-Margin Hyperplane

What if data are not linearly seperable?

— Still want to minimize $\|w\|$ (maximize margin)

— Would also like to minimize a loss function that penalizes points for being on the "wrong side"

# Hinge Loss Function

$$\text{hinge loss} = \sum_{i=1}^{n} \max \left[ 0, 1 - y_i (w^T x_i + b) \right]$$

$$= \begin{cases} 0 & \text{if } x \text{ outside or on margin} \\ > 0 & \text{if } x \text{ within margin} \end{cases}$$

Hinge loss penalizes misclassified points!

# Maximum Margin Hyperplane

Put "simply" want to minimize

$$C \times [\text{hinge loss}] + \left[ \frac{1}{\text{margin width}} \right]$$

where $C$ is a hyperparameter.

# Maximum Margin Hyperplane

Put "simply" want to minimize

$$[\text{hinge loss}] + \frac{1}{C}\left[\frac{1}{\text{margin width}}\right]$$

$$\sum_{i=1}^{N}\max\left(0, 1 - y_i(w^T x_i + b)\right) + \frac{1}{C}||w||^2$$

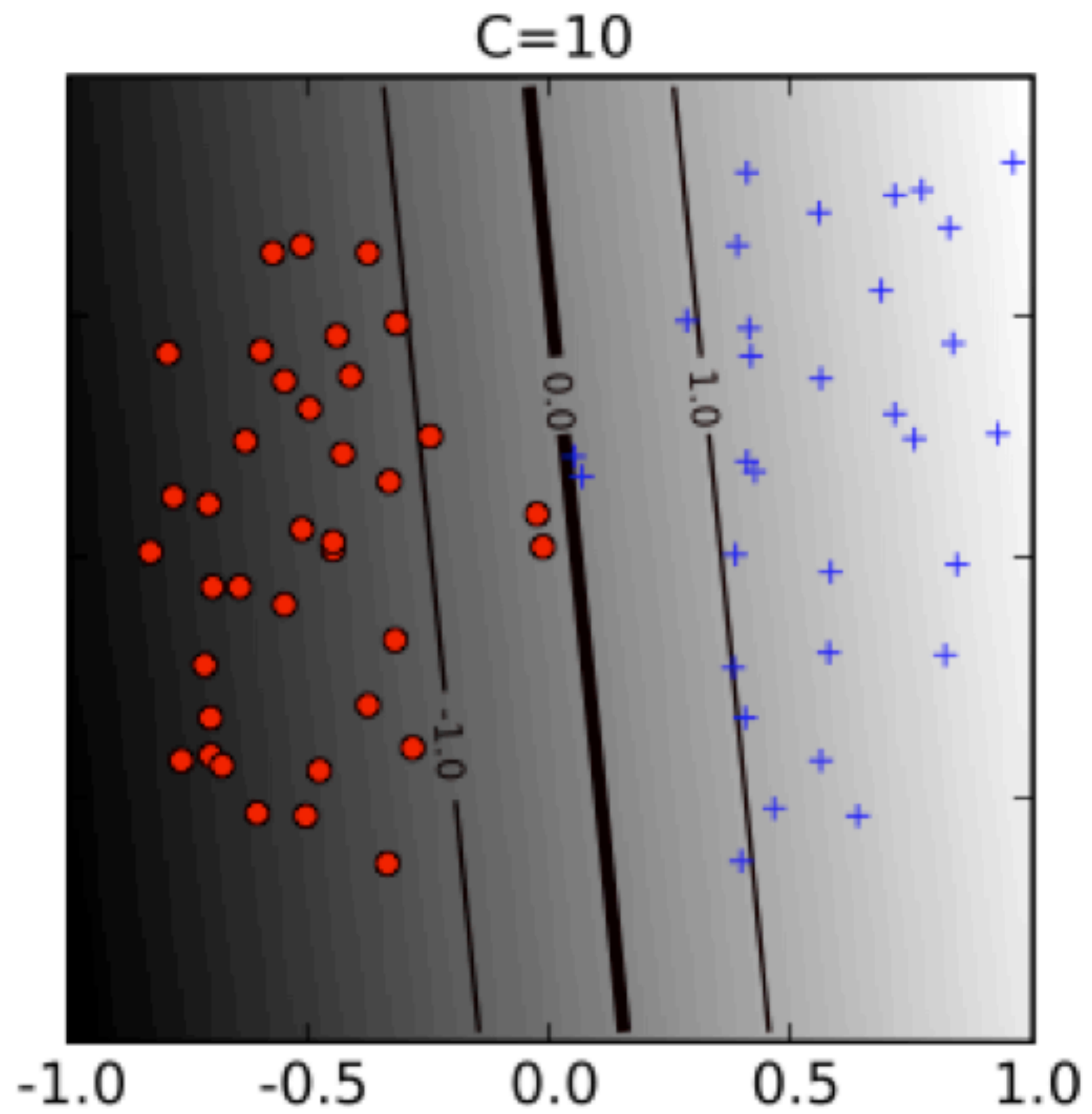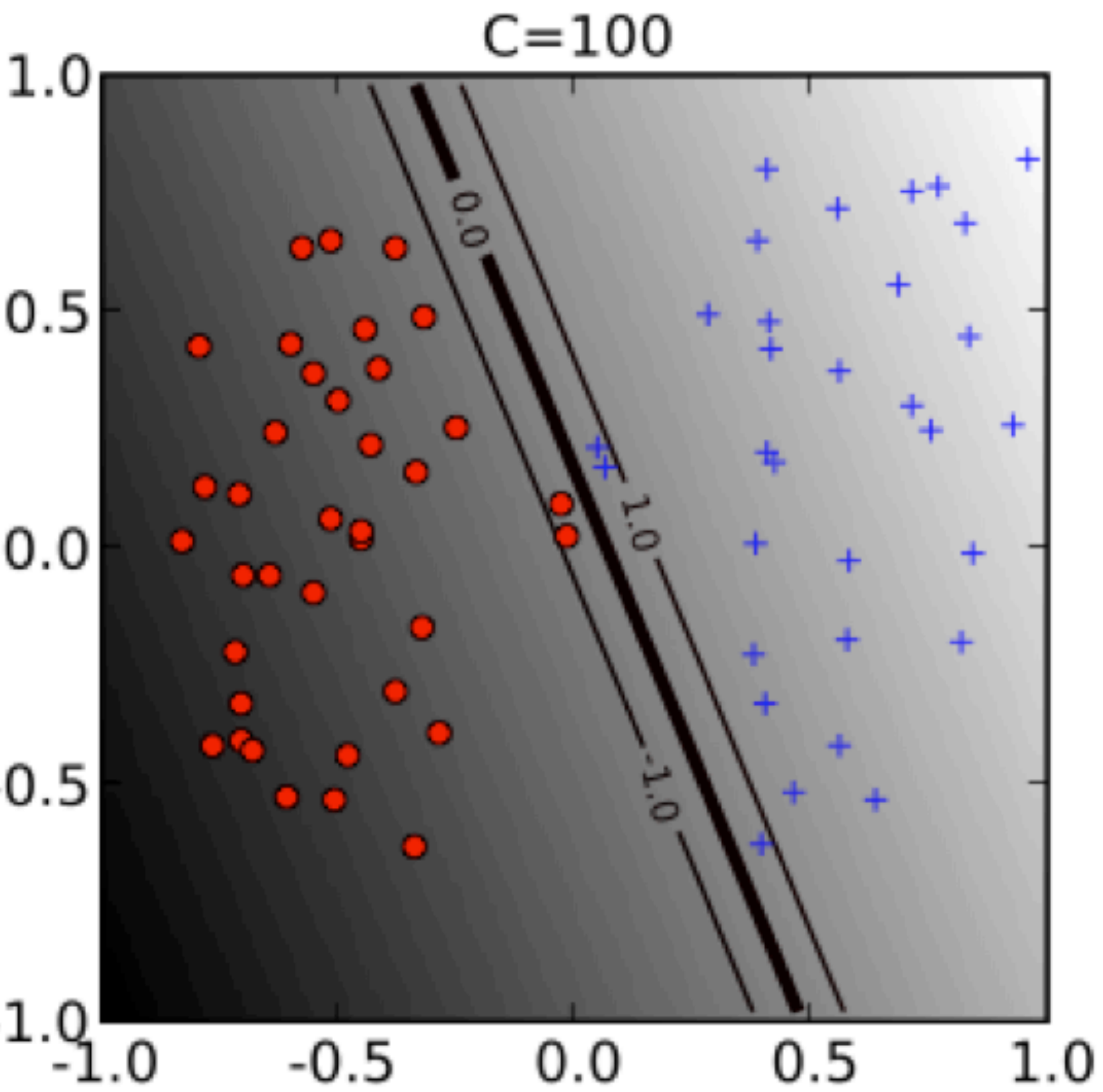i.e., loss + 1/C * regularization (c.f. Ridge!)

# Maximum Margin Hyperplane

Takeaway: **Bias/variance trade-off** is handled via the hyperparameter $C$

$$C \times [\text{hinge loss}] + \left[ \frac{1}{\text{margin width}} \right]$$

# Maximum Margin Hyperplane

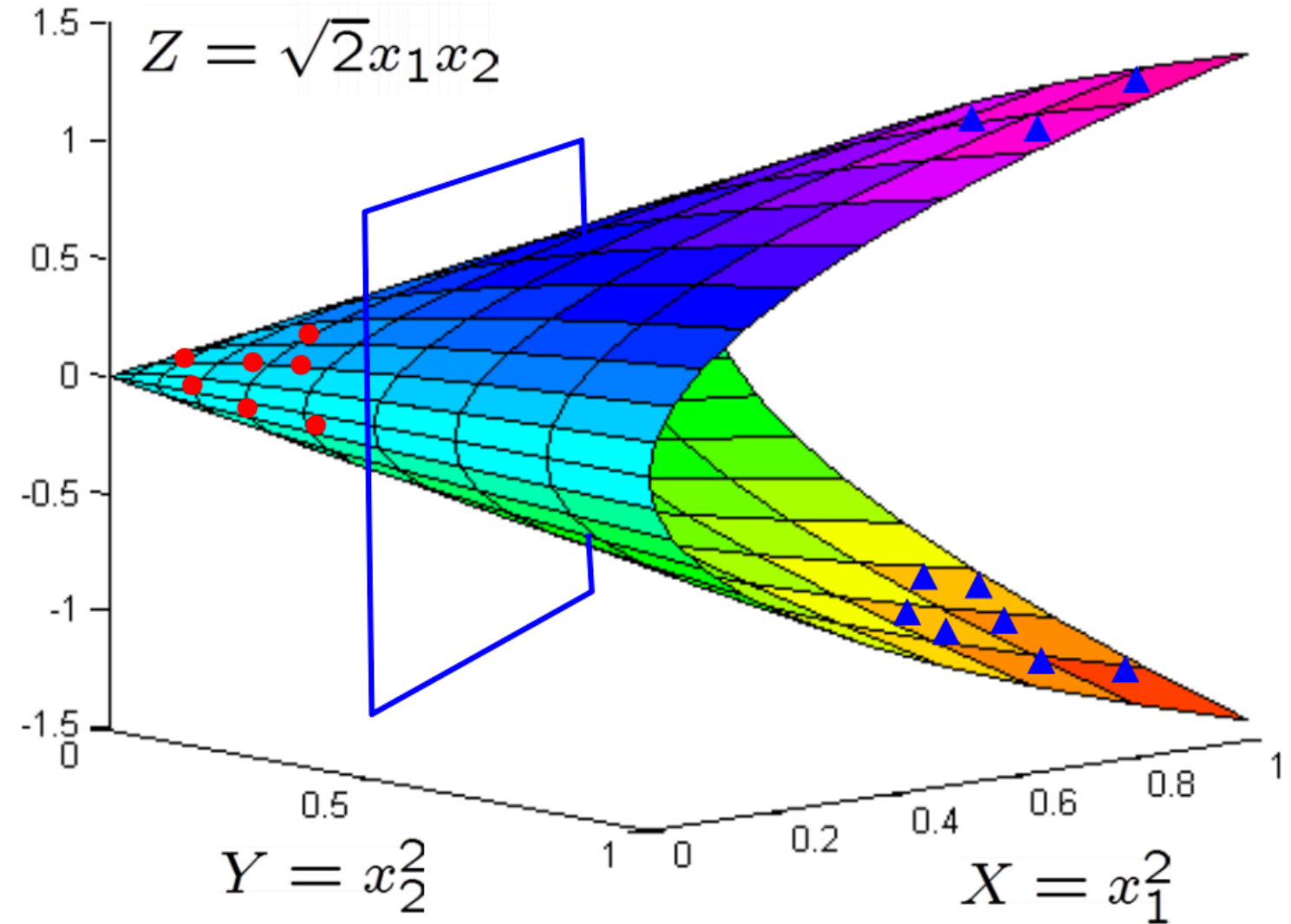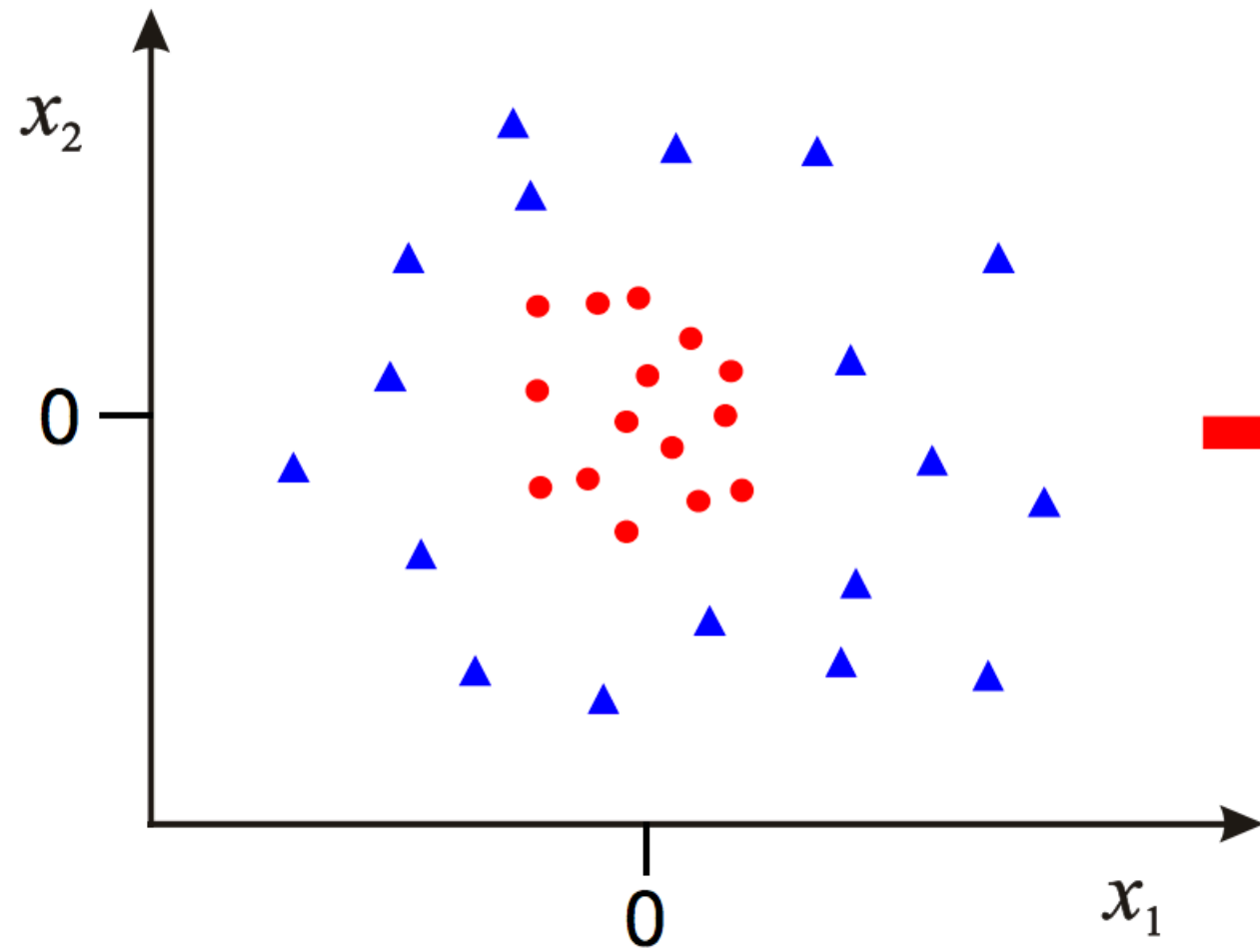$$C \times [\text{hinge loss}] + \left[ \frac{1}{\text{margin width}} \right]$$

— Large $C \rightarrow$ narrow margin, less tolerant of misclassification, tends toward high variance

— Small $C \rightarrow$ wider margin, more tolerant of misclassification, tends toward high bias

# What if your data are not separable?

Like, no where <span style="color:red">close</span> to linearly separable?

$$\Phi : \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \rightarrow \begin{pmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2}x_1x_2 \end{pmatrix} \qquad \mathbb{R}^2 \rightarrow \mathbb{R}^3$$



- Data is linearly separable in 3D
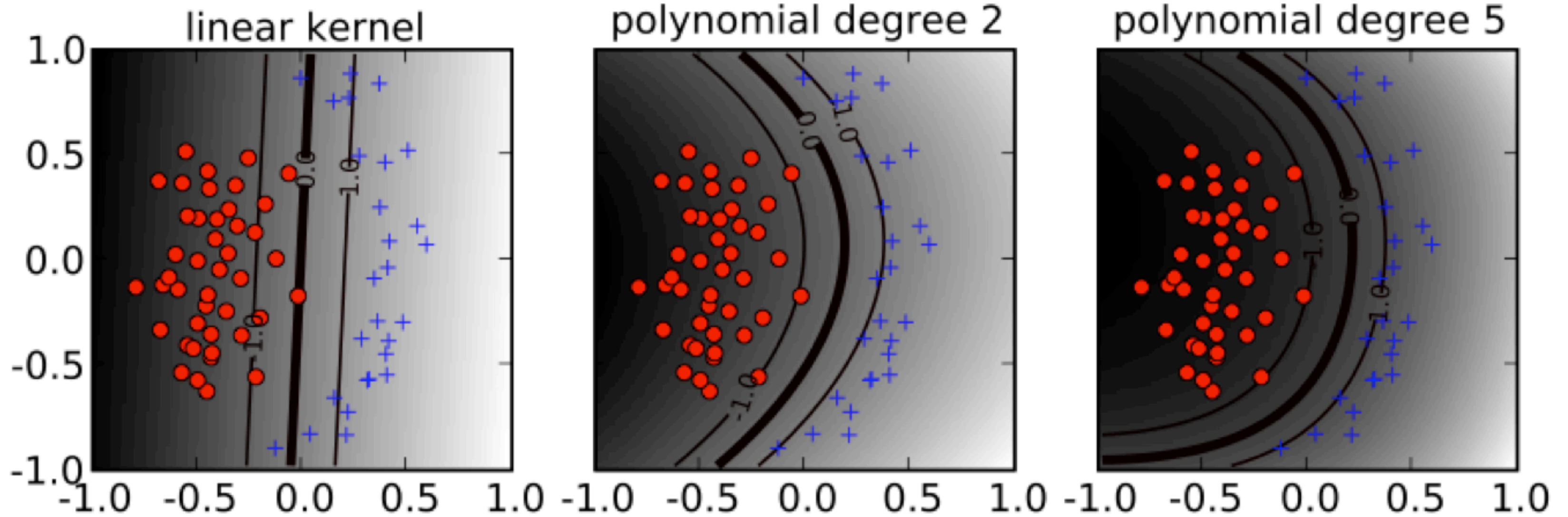- This means that the problem can still be solved by a linear classifier

21

# What if you data are not separable?
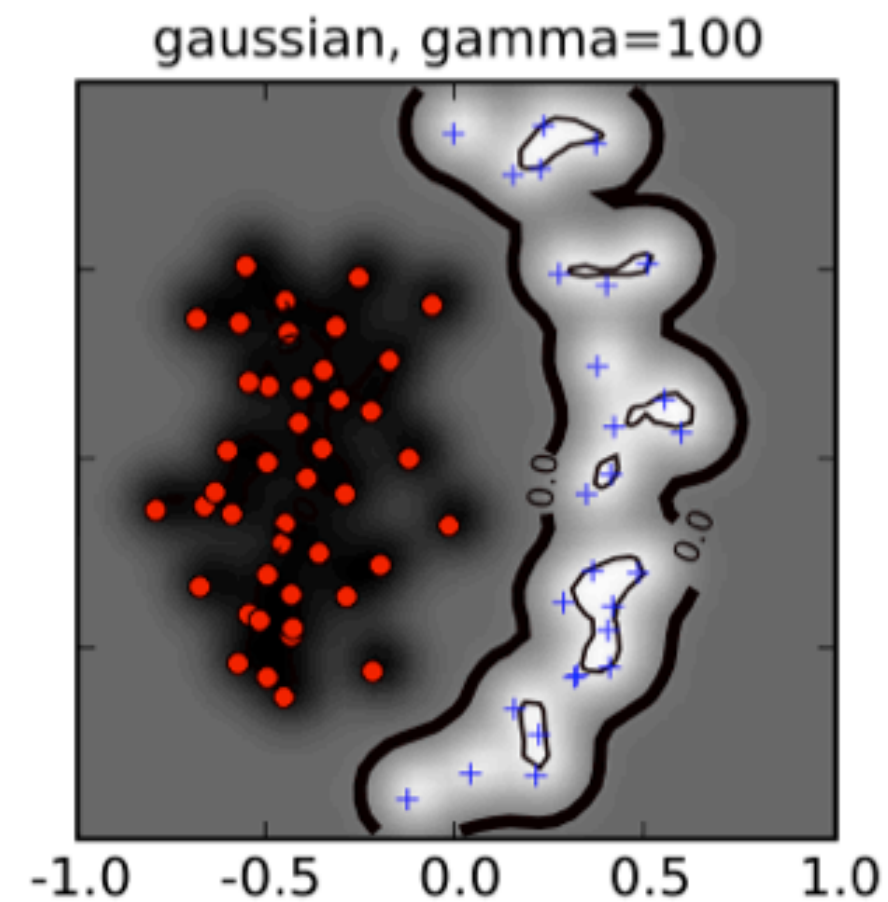
Kernel trick:

Replace

$$x \leftarrow \Phi(x)$$

$$w \leftarrow \Phi(w)$$
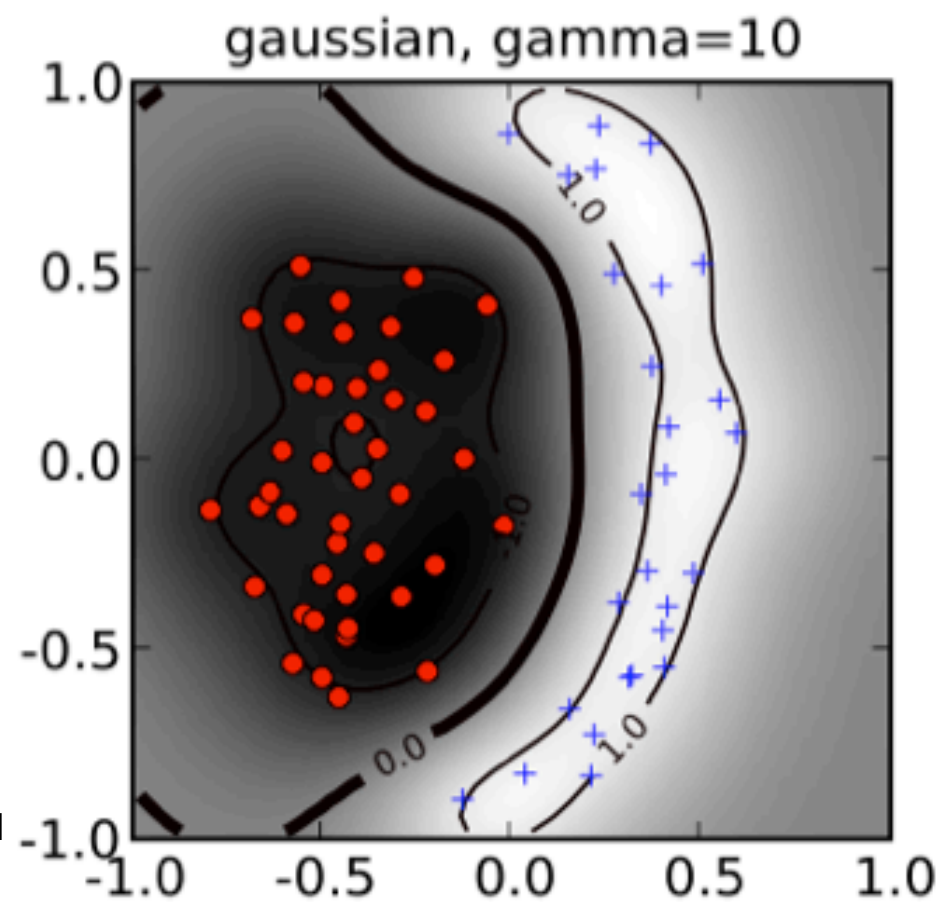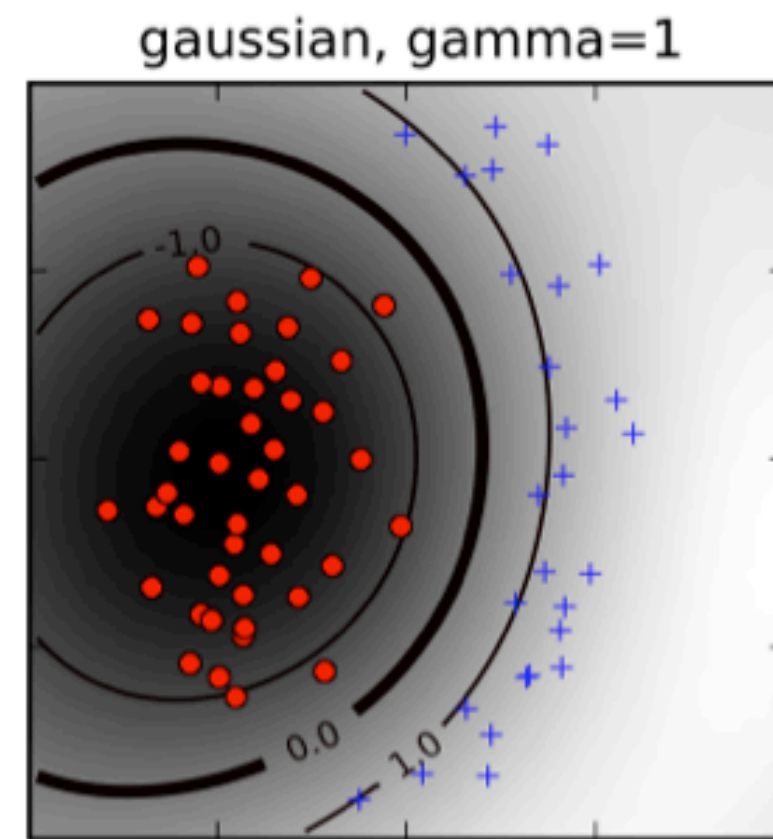
# Kernel SVM

— Linear

— Polynomial

— (Gaussian) Radial Bassis Function (RBF)

Four contour plots comparing Gaussian kernel SVM decision boundaries: "gaussian, gamma=0.1", "gaussian, gamma=1", "gaussian, gamma=10", and "gaussian, gamma=100".

# Pros and cons

## Pros
- Exceptional perfomance (historically widely used)
- Robust to outliers
- Effective in high dimensional data
- Can work with non-linearities
- Fast to compute even on non-linear (kernel trick)
- Low risk of overfitting

# Pros and cons

## Cons
- Blackbox
- Can be slow on large datasets

# When to use SVM vs. Logistic Regression

**Advice from Andrew Ng:**

**If there are more feature than training samples:**
 Use logistic regression or SVM without a kernel ("linear kernel")

**If there are about 10 times as many samples as features:**
 Use SVM with a Gaussian kernel

**If there are many more training samples than features:**
 Spend time feature engineering, then use logistic regression or SVM without a kernel

# Additional resources

**Taken from** `introduction-to-svm.ipynb` **in rep.**

— <u>For a really great resource check out these slides (some of which are cannabalized in this lecture).</u>

— <u>This website is also a great resource, on a slightly more technical level.</u>

— SVM docs on <u>SKLearn</u>

— Iris example on <u>SKLearn</u>

— Hyperplane walkthrough on <u>SKLearn</u>

— A comprehensive <u>user guide</u> to SVM. My fav!

# Additional resources

— A <u>blog post tutorial</u> of understanding the linear algebra behind SVM hyperplanes. Check <u>part 3</u> of this blog on finding the optimal hyperplane

— This <u>Quora discussion</u> includes a high-level overview plus a <u>20min video</u> walking through the core "need-to-knows"

— A <u>slideshow introduction</u> to the optimization considerations of SVM

# Additional resources

**Taken from** `introduction-to-svm.ipynb` **in rep.**

— Andrew Ng's <u>notes</u> on SVM from CS 229

— A <u>FULL LECTURE</u> (1hr+) from one of my fav lecturers (Dr Yasser) on SVM. He does a followup on <u>kernel tricks</u> too

— A <u>FULL LECTURE</u> (50min) (from MIT Opencoursewar)

— An infamous <u>paper</u> (cited 7000+ times!) on why SVM is a great text classifier

— An <u>advanced discussion</u> of SVMs as probabilistic