# Statistical Methods in AI (CSE 471)
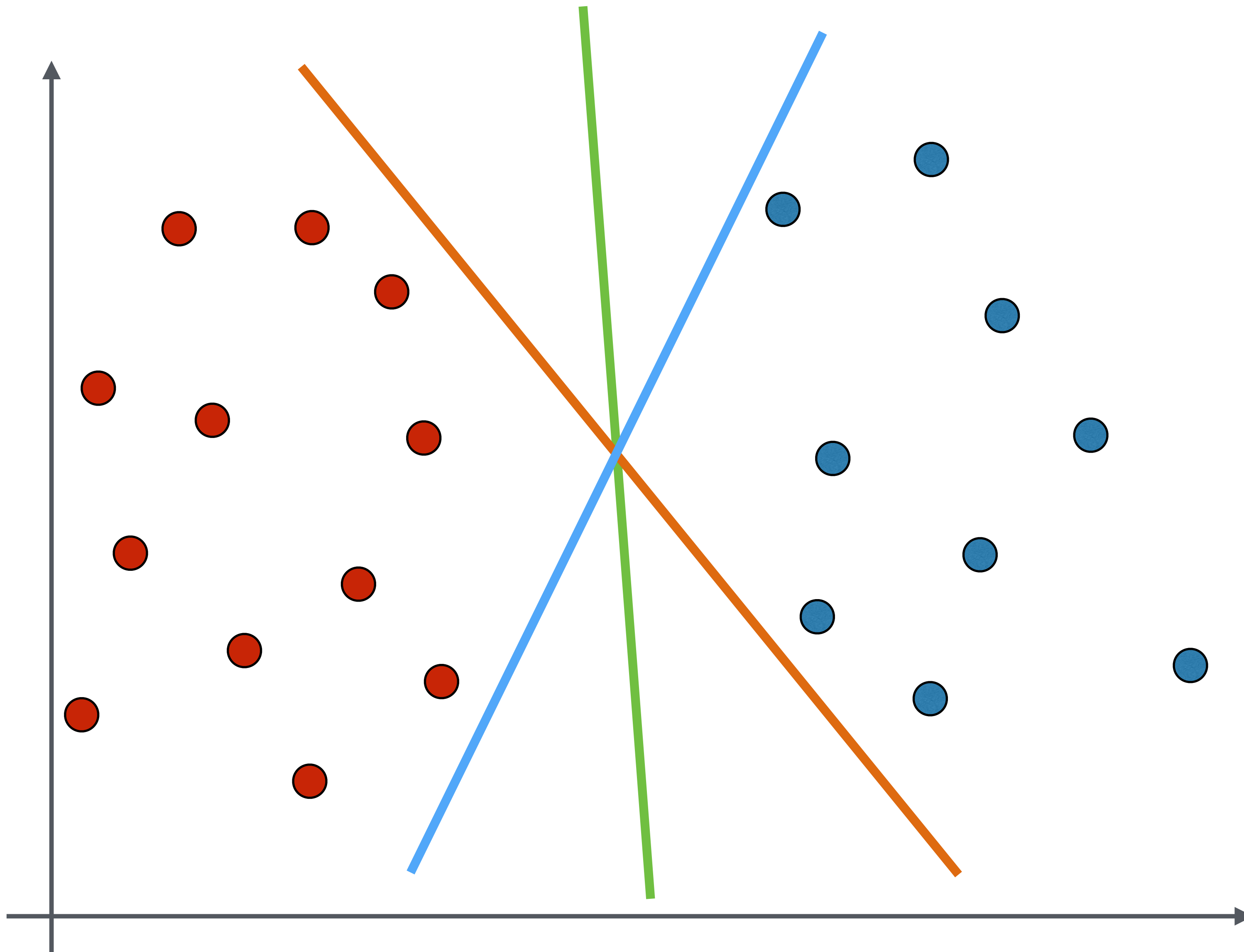# Lecture 6: Support Vector Machines
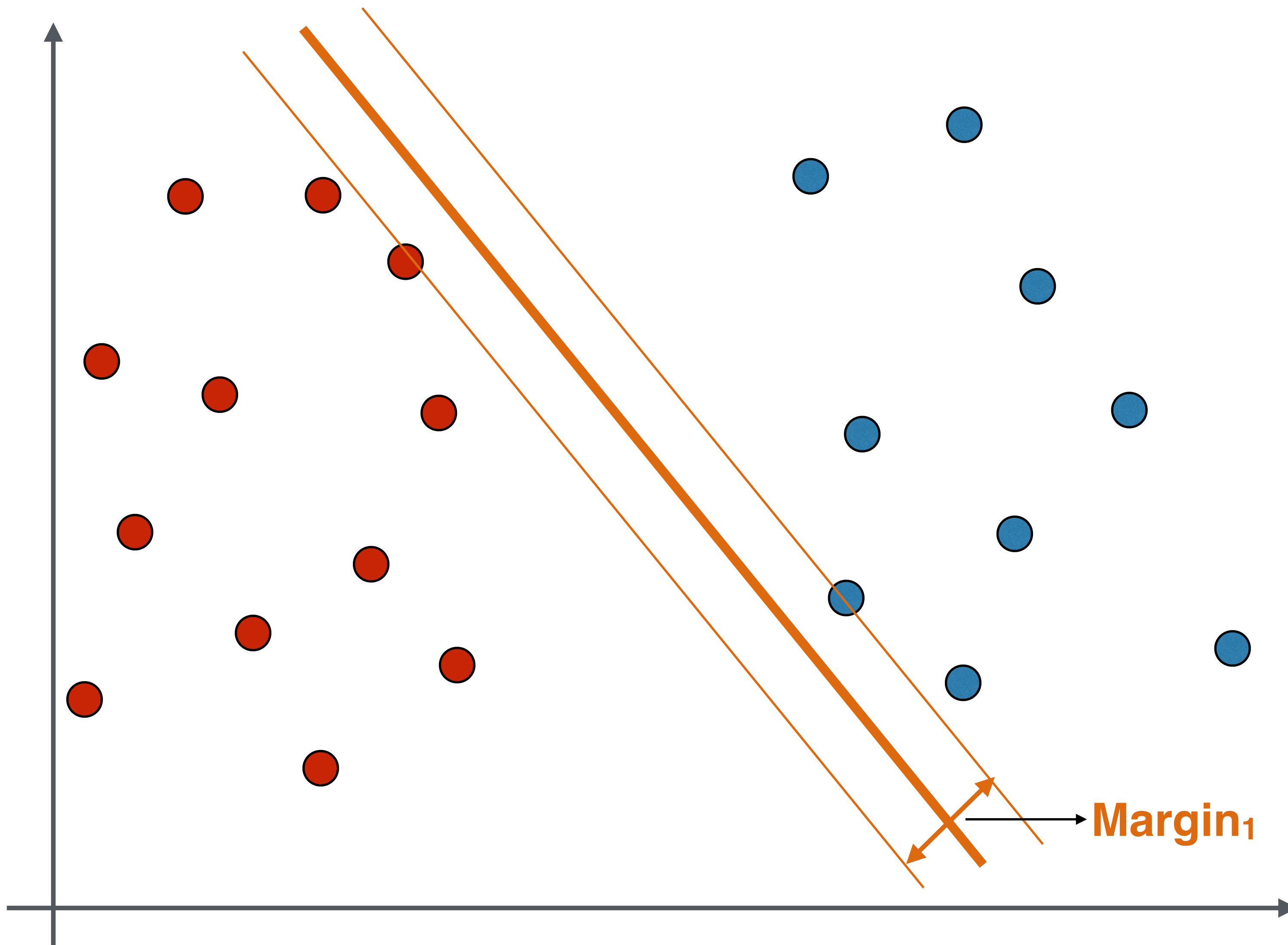
Vineet Gandhi
Centre for Visual Information Technology (CVIT)



INTERNATIONAL INSTITUTE OF
INFORMATION TECHNOLOGY
H Y D E R A B A D

# Multiple solutions exist for linearly separable data



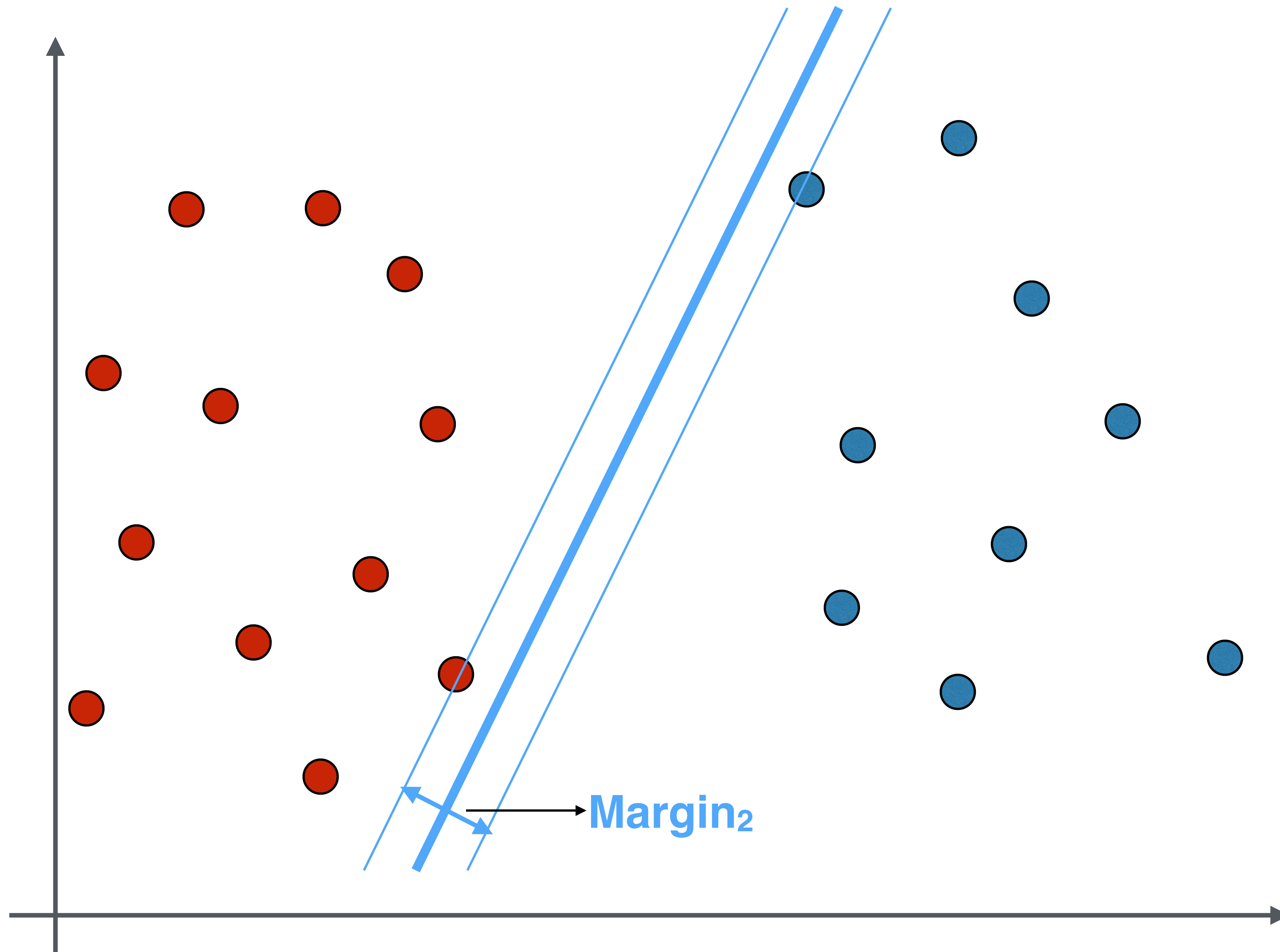Are all solutions equally good?
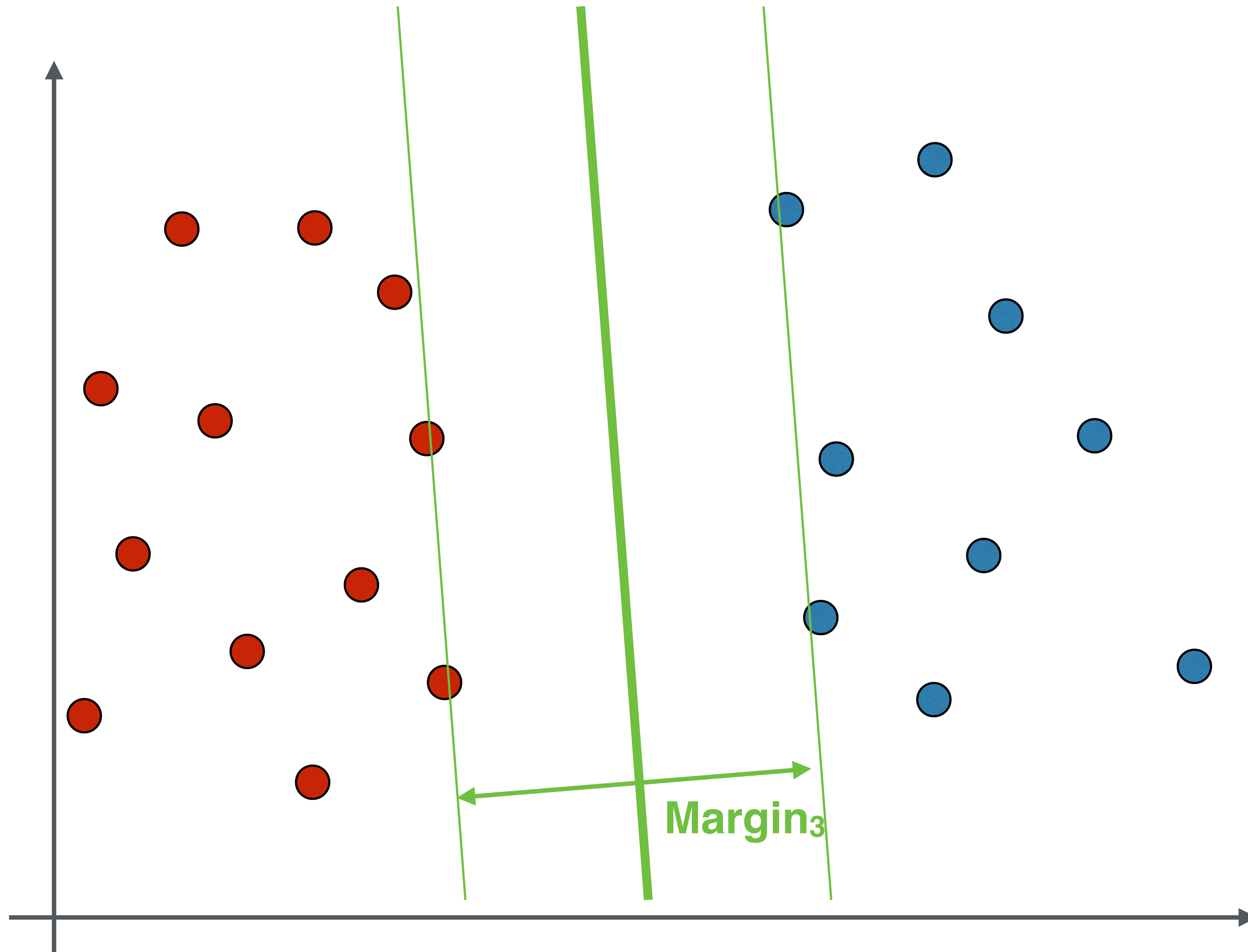
# Margin: No-mans Band



Margin: Width of a band around decision boundary without any training samples

**Margin₁**

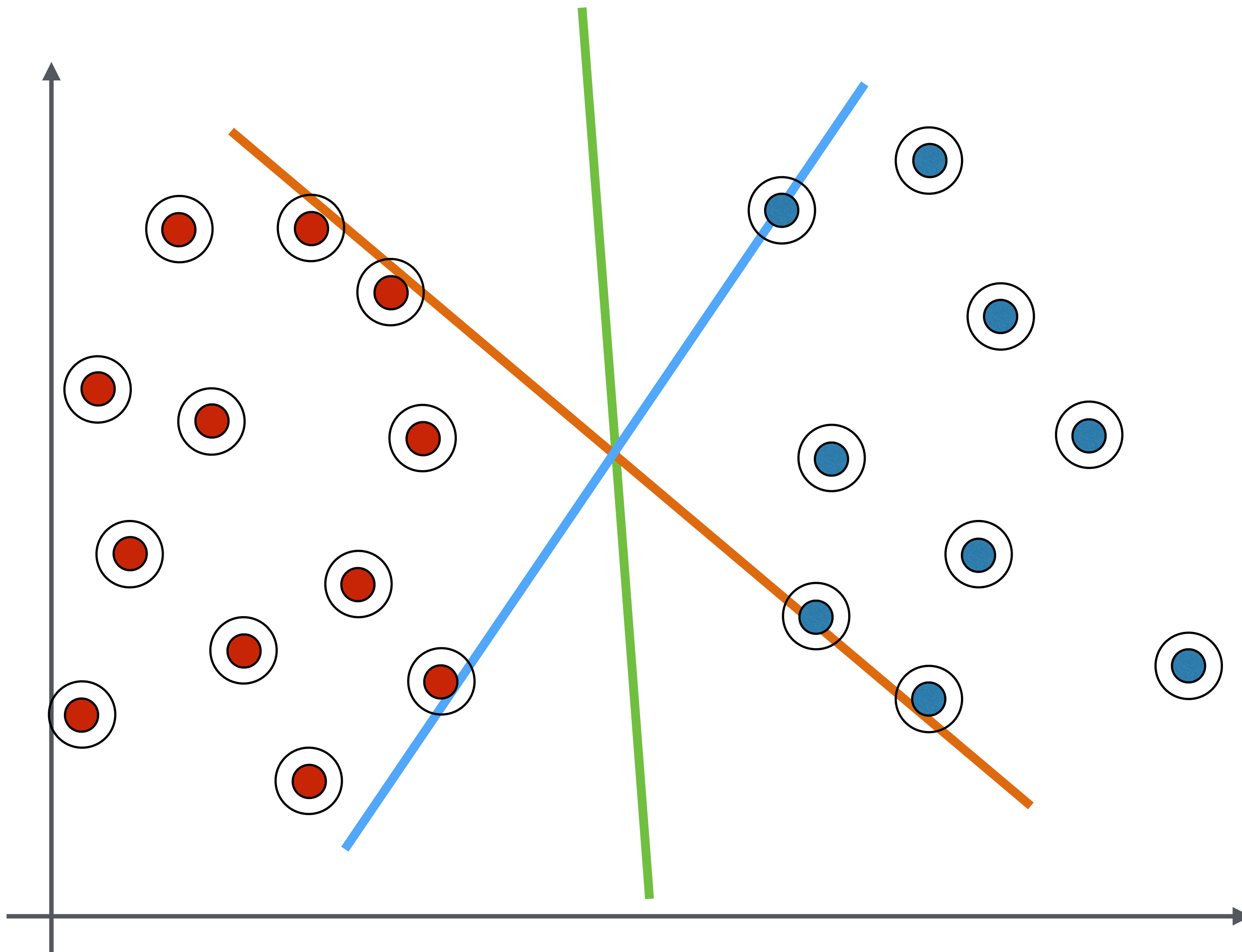# Multiple solutions exist for linearly separable data



Margin: Width of a band around decision boundary without any training samples

**Margin₂**
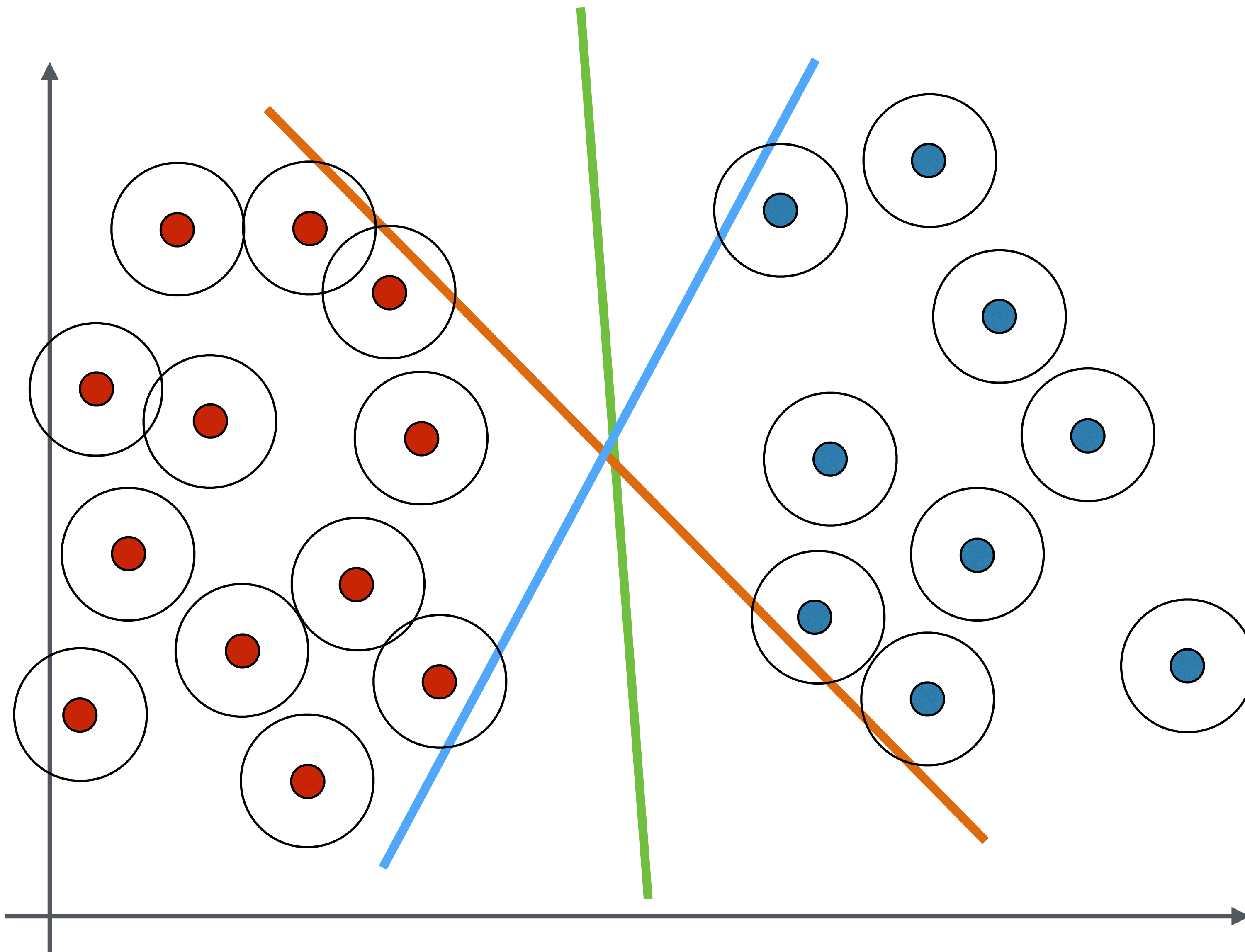
# Multiple solutions exist for linearly separable data



Is a Larger Margin better? Why?

Margin$_3$

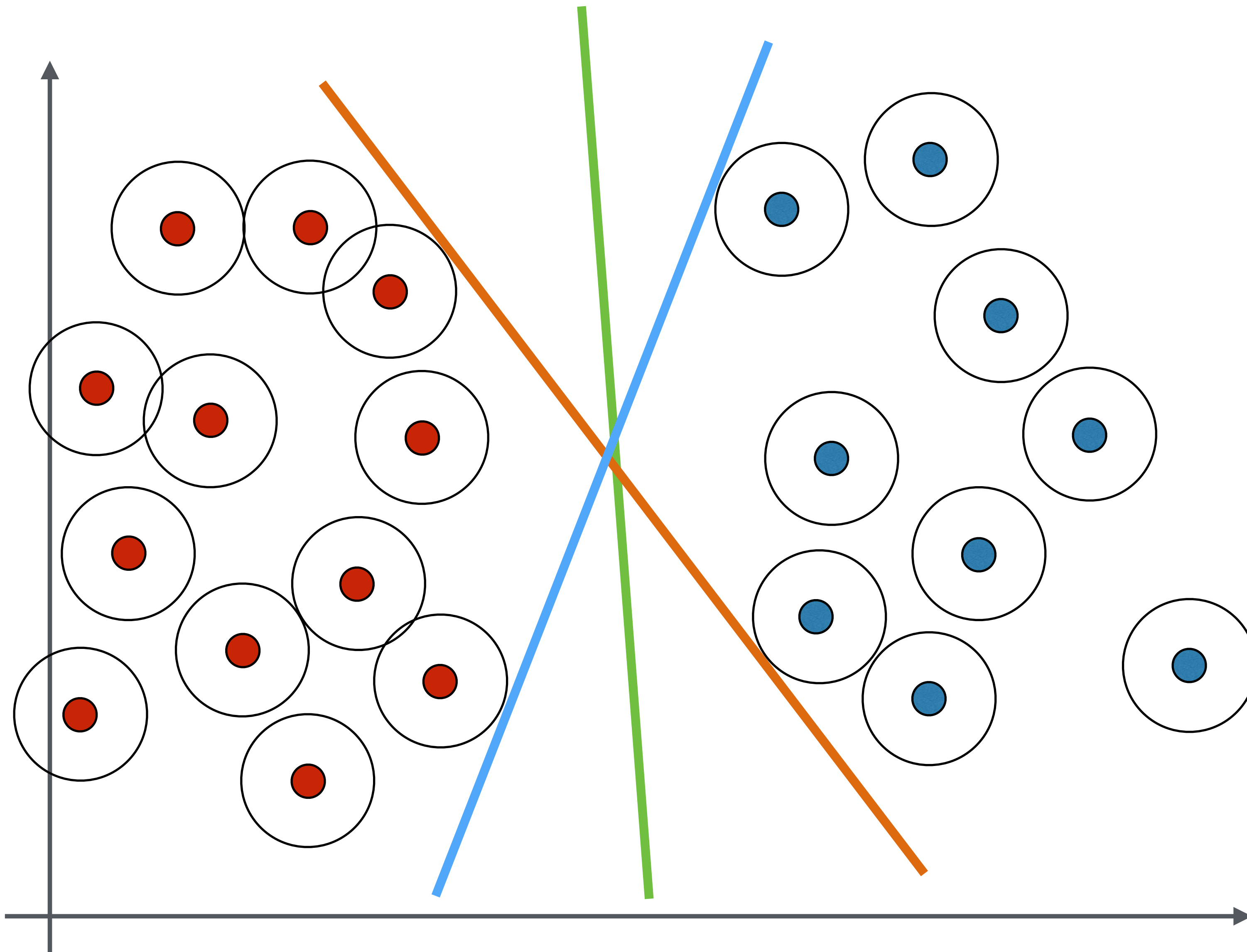# Margin: Bubbles around samples



Margin: Radius of a region around each training sample, through which the decision boundary cannot pass

# Margin: Bubbles around samples



Margin: Radius of a region around each training sample, through which the decision boundary cannot pass

# Margin: Bubbles around samples



Margin: Radius of a region around each training sample, through which the decision boundary cannot pass

As the margin increases, the feasible region reduces

# Margin: Bubbles around samples



Margin: Radius of a region around each training sample, through which the decision boundary cannot pass

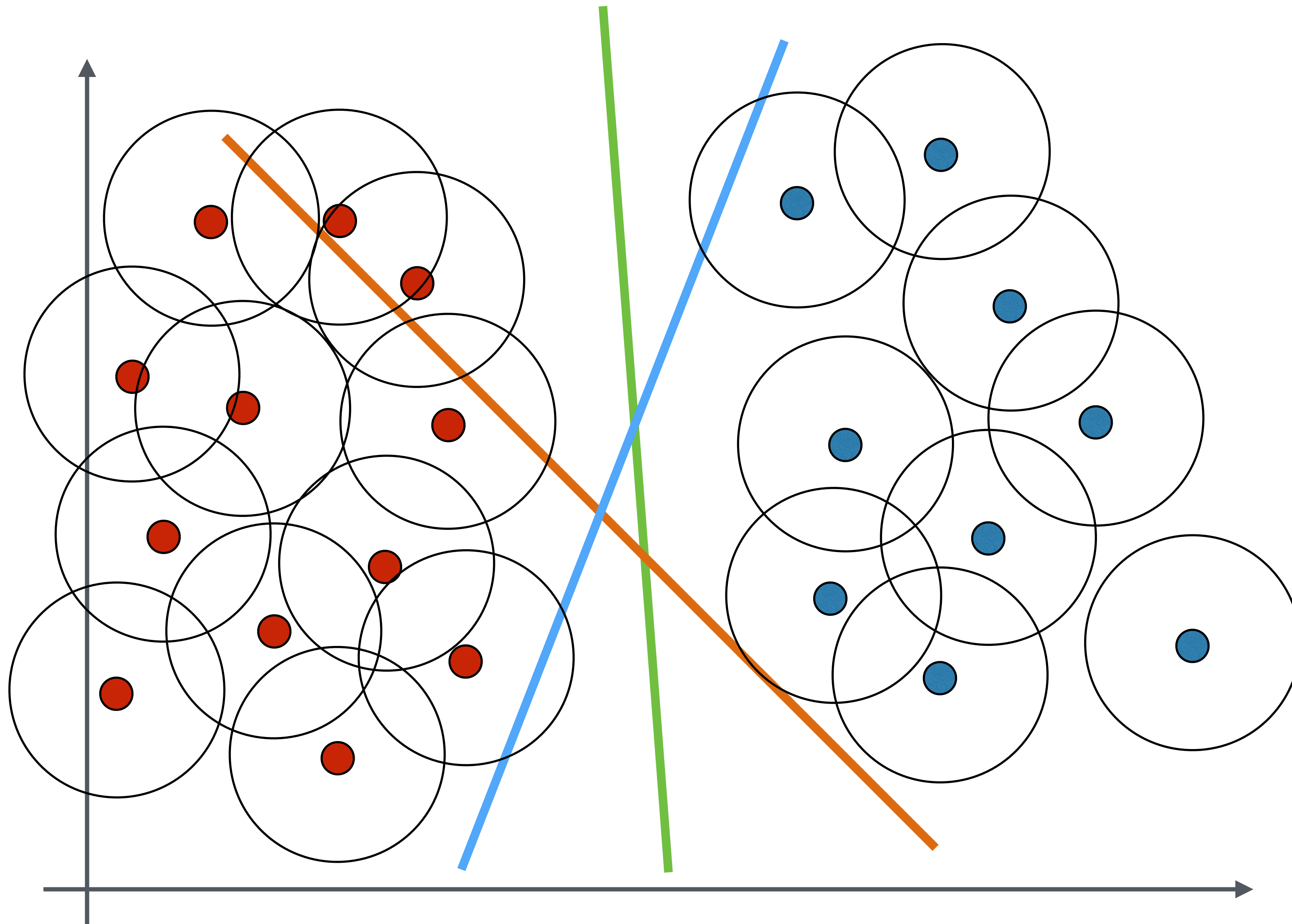As the margin increases, the feasible region reduces

# Margin: Bubbles around samples



Margin: Radius of a region around each training sample, through which the decision boundary cannot pass

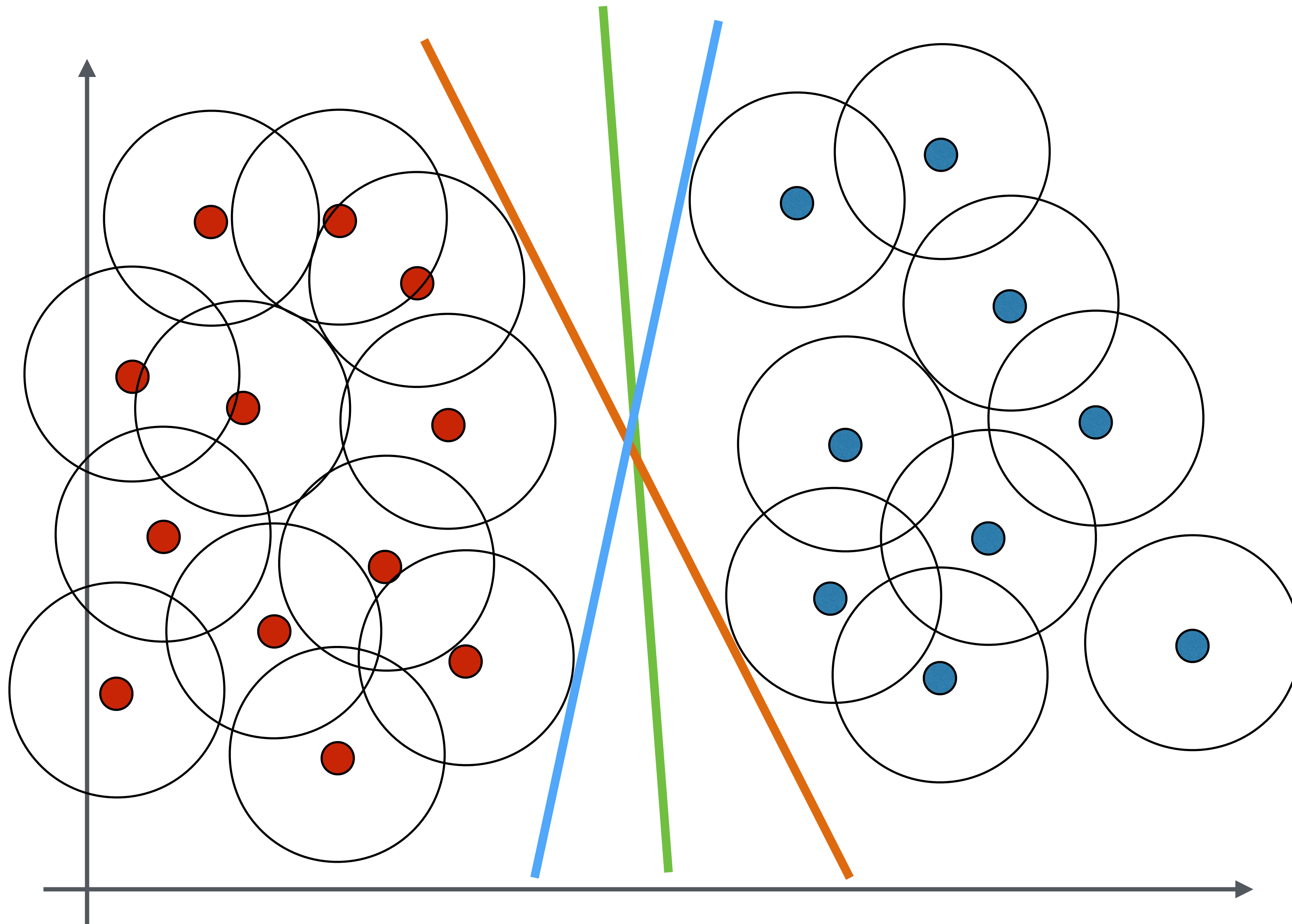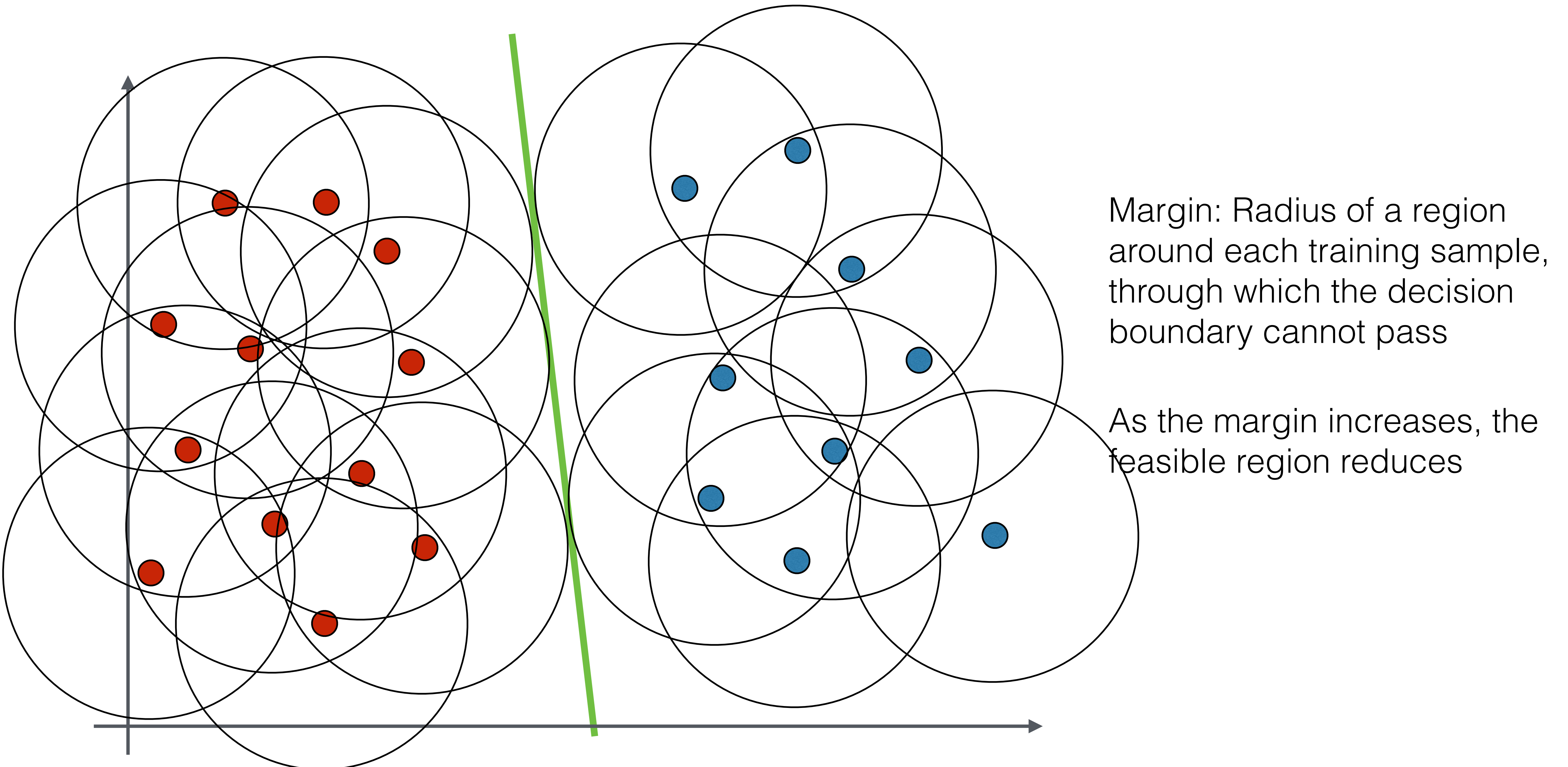As the margin increases, the feasible region reduces

# Margin: Bubbles around samples

Margin: Radius of a region around each training sample, through which the decision boundary cannot pass

As the margin increases, the feasible region reduces

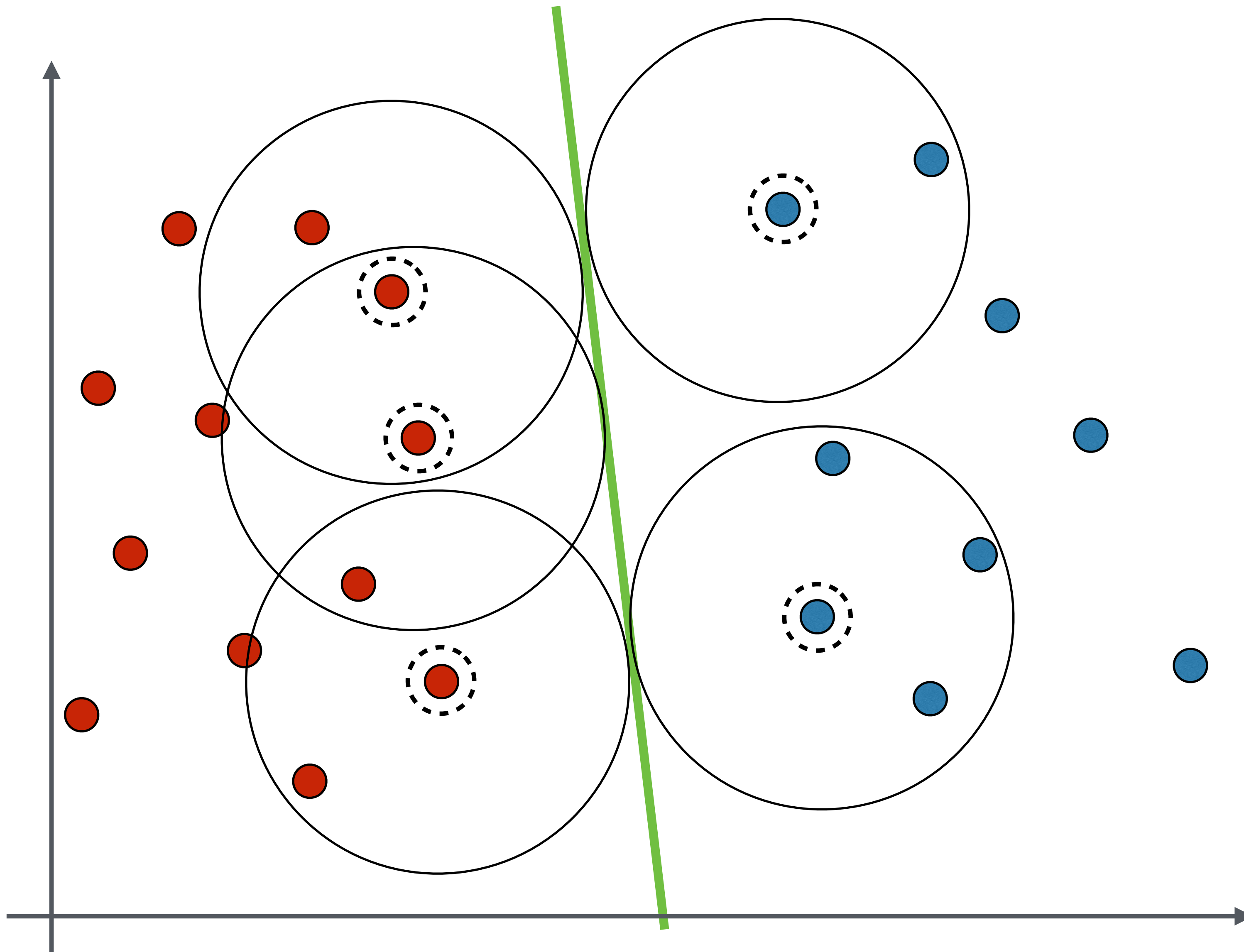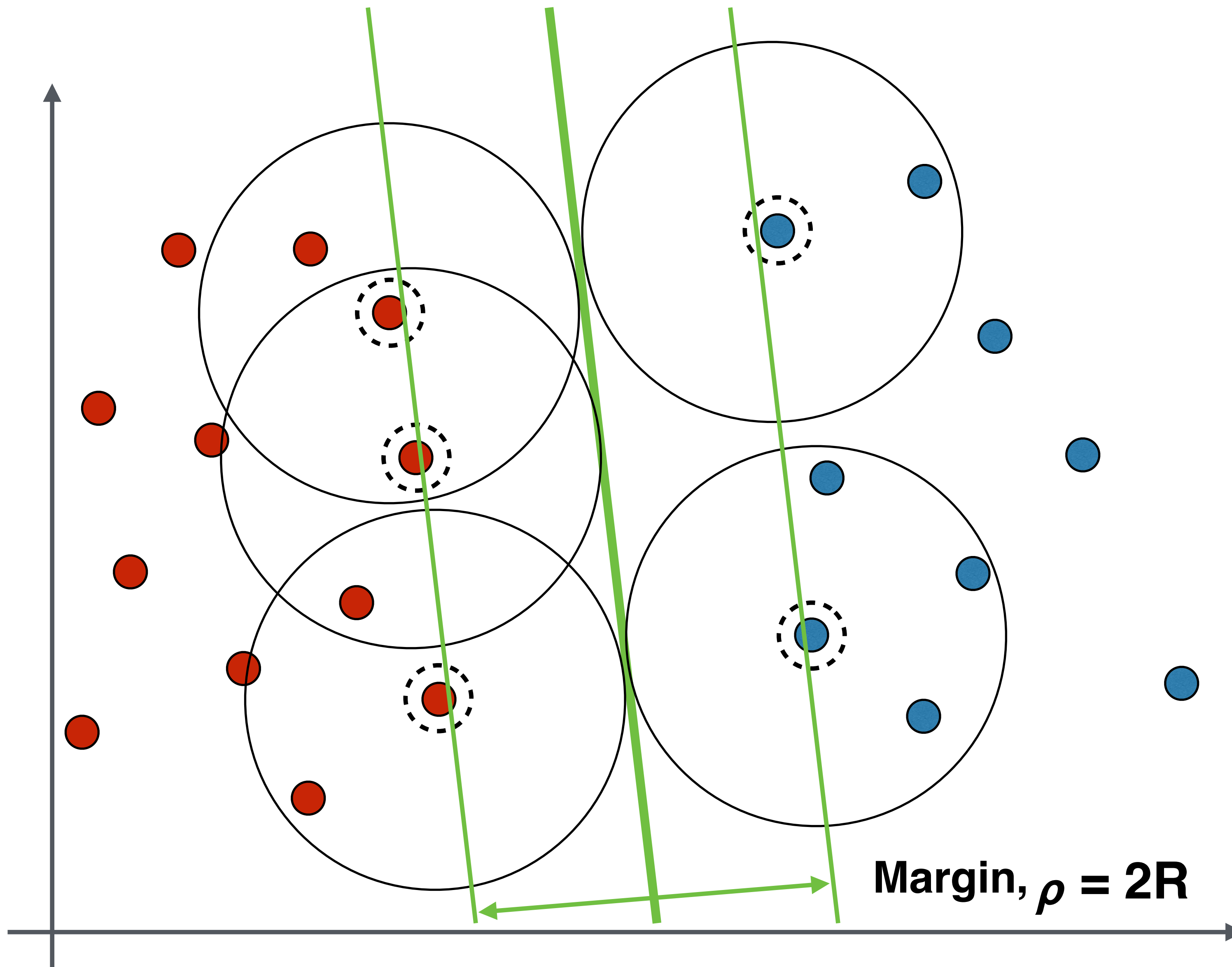# Margin: Bubbles around samples
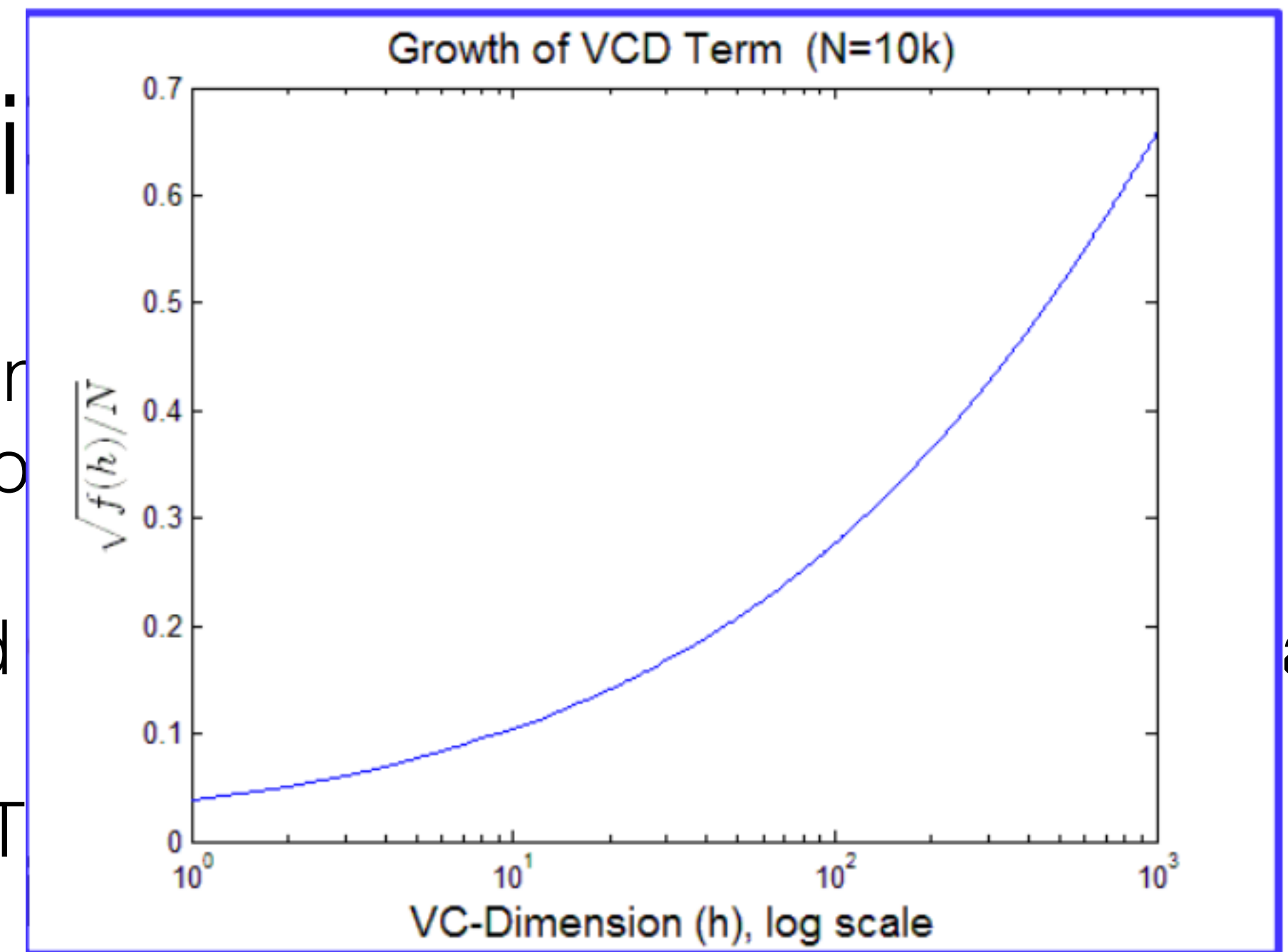


A few samples control the Decision Boundary

# Band vs. Bubbles



Samples that support the boundary are called **_Support Vectors_**

Both interpretations lead yield the same decision boundary

Margin, $\rho$ = 2R

# Break though work from Vapni[k]



Growth of VCD Term (N=10k)

1. Vapnik, Vladimir N., and A. Ya Chervonenkis. "On the u[niform convergence of relative] frequencies of events to their probabilities." Measures o[f ... ] Primenen., 1971, Volume 16, Issue 2, Pages 264–279
2. Vapnik, Vladimir N., Estimation of Dependences Based [...a,] Moscow.
3. Vapnik, Vladimir N., The Nature of Statistical Learning T[...]

**Bound on expected loss:**
$$R(\alpha) \le R_{train}(\alpha) + \sqrt{f(h) \Big/ N}$$

$h$ is the VC dimension, and $f(h)$ is given by:

$$f(h) = h + h\log(2N) - h\log(h) - c$$

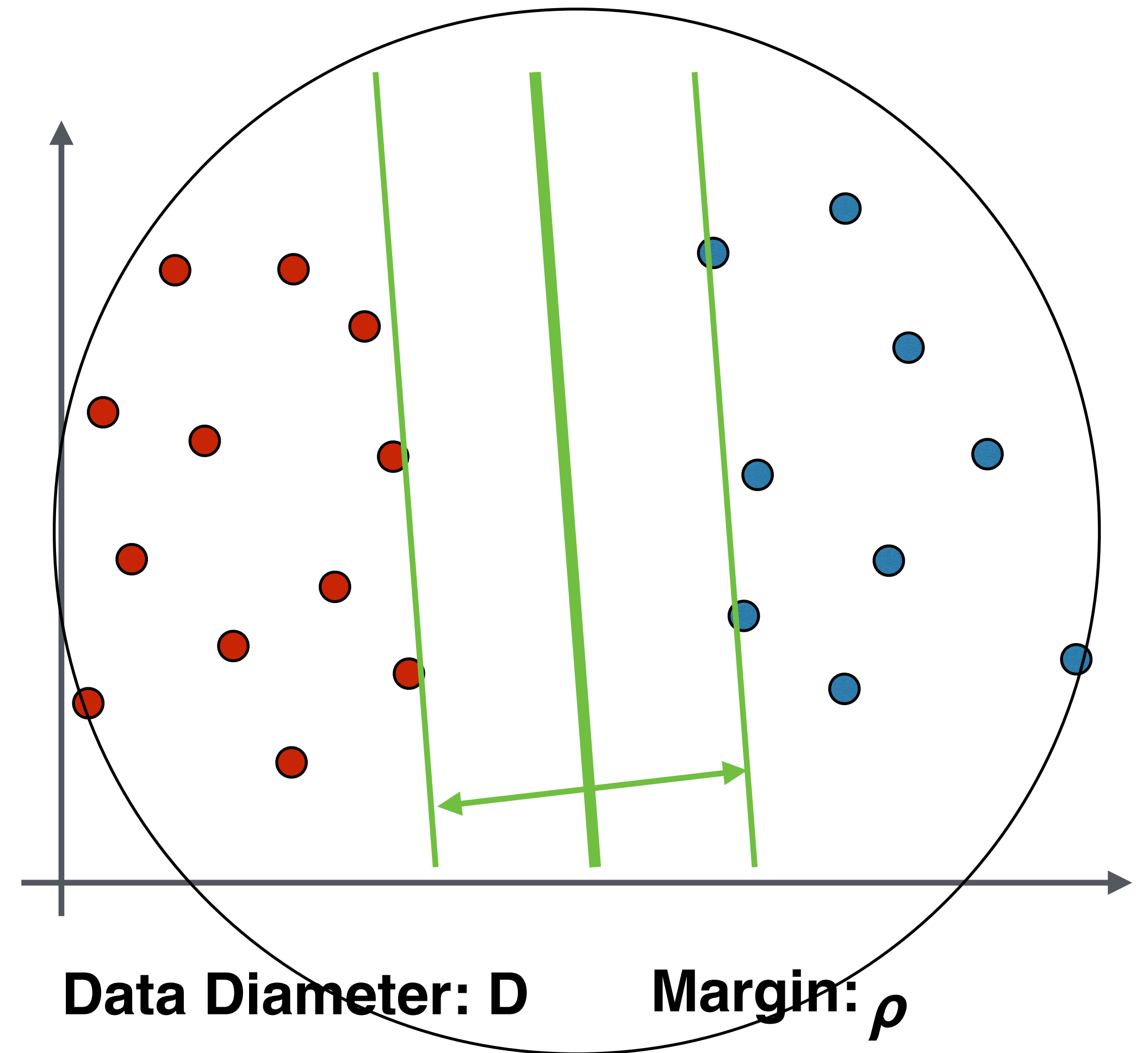# Why maximise the margin?

- To reduce test error, keep training error low (say 0), and minimize the VC-dimension, $h$.

Relative Margin : $\rho\!\!\big/_{\!\!D}$

VC-D, $h \le \min\left\{ d \;,\; \left\lceil \dfrac{D^2}{\rho^2} \right\rceil \right\} + 1$

- Maximizing margin improves generalization.

- $h$ can be made independent of the dimensionality: $d$.
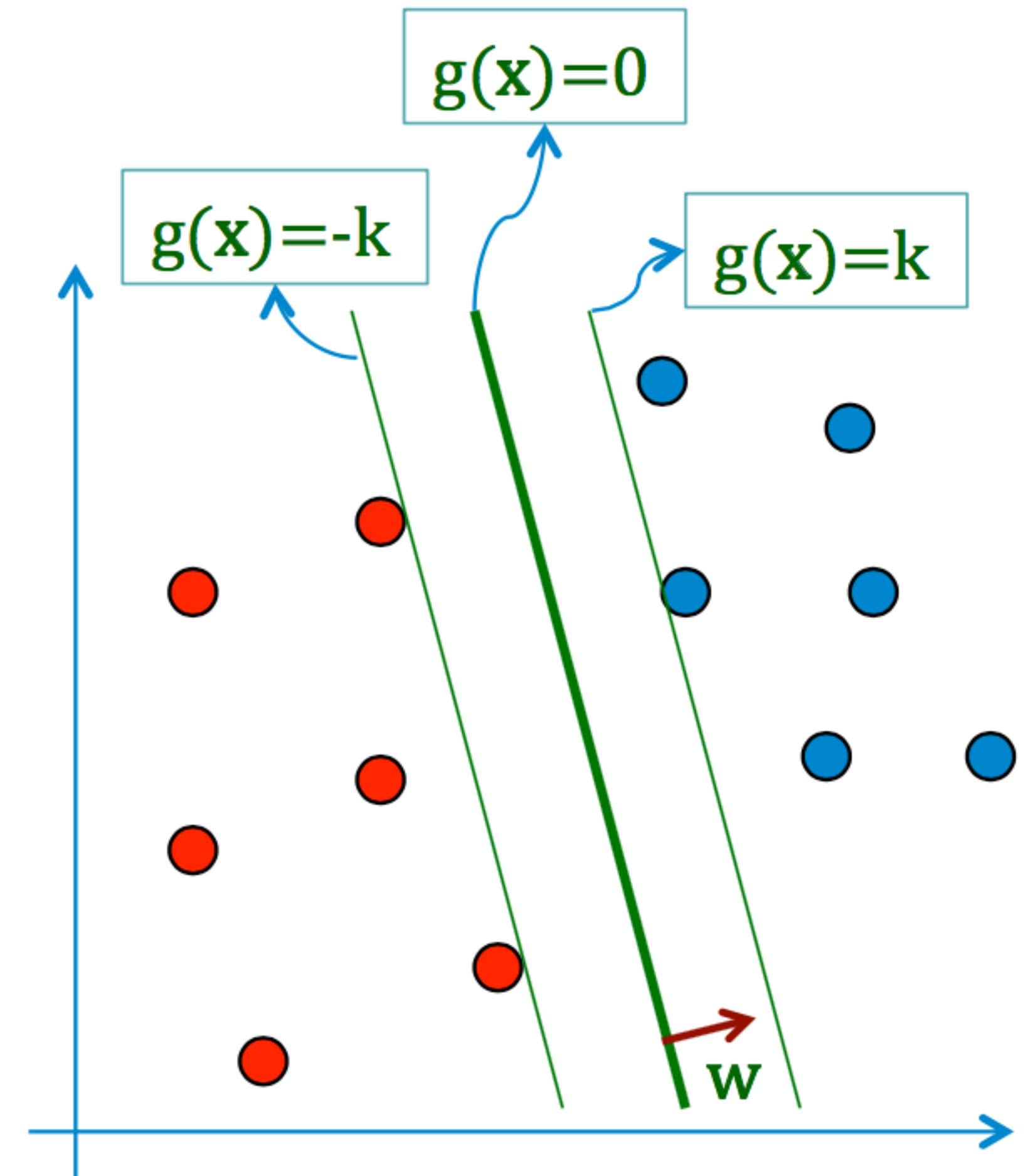


**Data Diameter: D**     **Margin:** $\rho$

# Formalizing the margin



Dec. Boundary: $W^T X + b = 0$

Note: The value of $W^T X_i + b$ is dependent on the scale of X and W

Let parallel hyperplanes be:
$W^T X + b = \pm\varepsilon$

# Formulation

- Let $g(x)=w^{T}x+b$.

- We want to maximize k such that:
  - $w^{T}x_{i}+b \geq k$ for $d_{i}=1$
  - $w^{T}x_{i}+b \leq -k$ for $d_{i}=-1$

- Value of $g(x)$ dependents on $\|w\|$ :
  1. Keep $\|w\|=1$, and maximize $g(x)$, or
  2. Let $g(x) \geq 1$, and minimize $\|w\|$.

- We use approach (2) and formulate the problem as:
  - Minmize: $\frac{1}{2}w^{T}w$
  - Subject to: $d_{i}(w^{T}x_{i}+b) \geq 1$, for $i=1..N$

# Optimization


Saddle Point

Minimize : $\Phi(\mathbf{w}) = \frac{1}{2}\mathbf{w}^T\mathbf{w}$

Subject to : $d_i(\mathbf{w}^T\mathbf{x_i} + b) - 1 \geq 0 \quad \forall i$
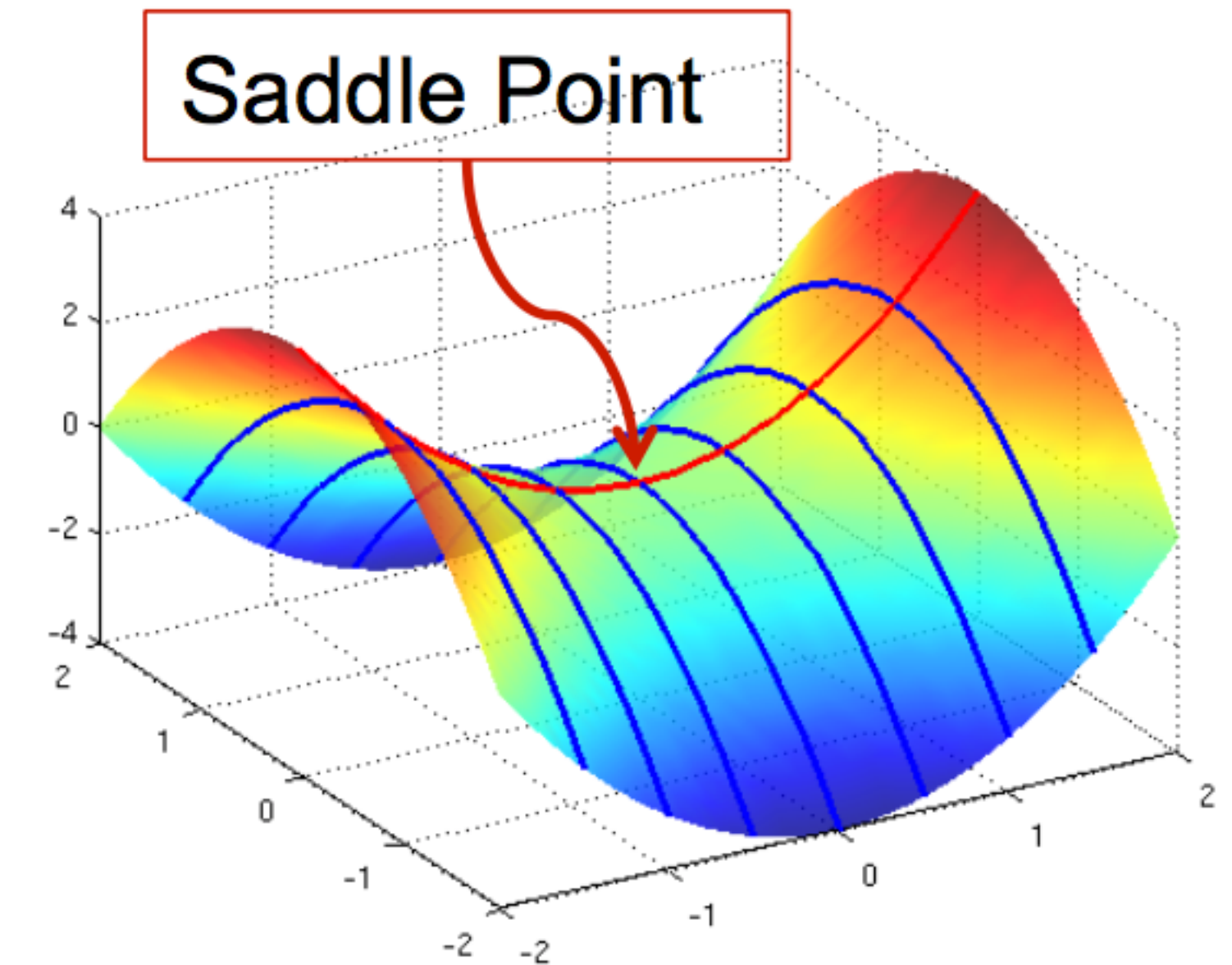
Quadratic function: QP solvers

Lagrangian form:

Minimize : $J(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2}\mathbf{w}^T\mathbf{w} - \sum_{i=1}^{N}\alpha_i d_i(\mathbf{w}^T\mathbf{x_i} + b) + \sum_{i=1}^{N}\alpha_i$

Subject to : $\alpha_i \geq 0 \quad \forall i$

Minimize J with respect to w and b, and maximize with respect to $\alpha$.

# Converting to Dual form

Objective: $J(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2}\mathbf{w}^T\mathbf{w} - \sum_{i=1}^{N}\alpha_i d_i(\mathbf{w}^T\mathbf{x_i} + b) + \sum_{i=1}^{N}\alpha_i$

At the optimum:

$1: \dfrac{\partial J}{\partial \mathbf{w}} = 0$ and $2: \dfrac{\partial J}{\partial b} = 0$

$1: \mathbf{w}_o = \sum_{i=1}^{N}\alpha_i d_i \mathbf{x_i}$ $\quad$ $2: \sum_{i=1}^{N}\alpha_i d_i = 0$ $\quad$ $3: \alpha_i[d_i(\mathbf{w}_o^T\mathbf{x_i} + b_o) - 1] = 0$

Obj: $J(\mathbf{w}, b, \boldsymbol{\alpha}) = \sum_{i=1}^{N}\alpha_i + \frac{1}{2}\mathbf{w}^T\mathbf{w} - \mathbf{w}^T\sum_{i=1}^{N}\alpha_i d_i\mathbf{x_i} - b\sum_{i=1}^{N}\alpha_i d_i$

$$Q(\boldsymbol{\alpha}) = \sum_{i=1}^{N}\alpha_i - \frac{1}{2}\sum_{i=1}^{N}\sum_{j=1}^{N}\alpha_i\alpha_j d_i d_j \mathbf{x_i}^T\mathbf{x}_j$$

# Solving the Dual form

$$Q(\boldsymbol{\alpha}) = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j d_i d_j \mathbf{x_i}^T \mathbf{x}_j$$

$$\text{Subject to } \alpha_i \geq 0 \quad \forall_i \quad \text{and} \quad \sum_{i=1}^{N} \alpha_i d_i = 0$$

QP Solver

$\alpha_i$

$$\mathbf{w}_o = \sum_{i=1}^{N} \alpha_i d_i \mathbf{x_i} \longrightarrow \alpha_i [d_i (\mathbf{w}_o^T \mathbf{x_i} + b_o) - 1] = 0 \longrightarrow b_o = 1 - \mathbf{w}_o^T \mathbf{x}_{s+}$$