# Support Vector Machines & Kernels

## Doing *really* well with linear decision surfaces

These slides were assembled by Eric Eaton, with grateful acknowledgement of the many others who made their course materials freely available online. Feel free to reuse or adapt these slides for your own academic purposes, provided that you include proper attribution. Please send comments and corrections to Eric.

# Outline

- Prediction
  - Why might predictions be wrong?
- Support vector machines
  - Doing really well with linear models
- Kernels
  - Making the non-linear linear
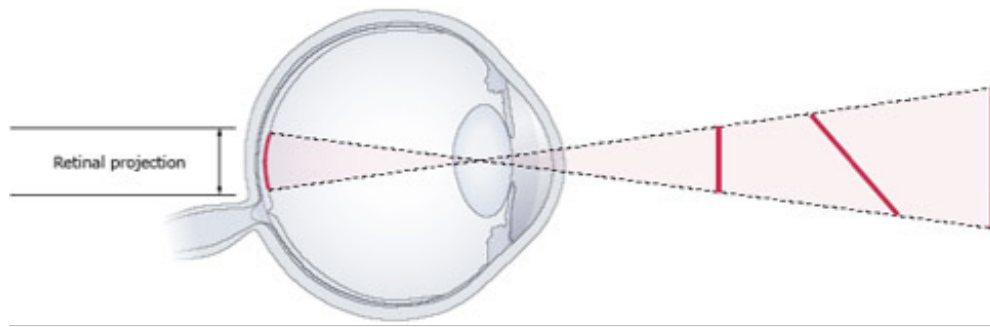
# Why Might Predictions be Wrong?

- True non-determinism
  - Flip a biased coin
  - $p(\text{heads}) = \boldsymbol{\theta}$
  - Estimate $\boldsymbol{\theta}$
  - If $\boldsymbol{\theta} > 0.5$ predict 'heads', else 'tails'

Lots of ML research on problems like this:
  - Learn a model
  - Do the best you can in expectation

# Why Might Predictions be Wrong?

- Partial observability
  - Something needed to predict $y$ is missing from observation x
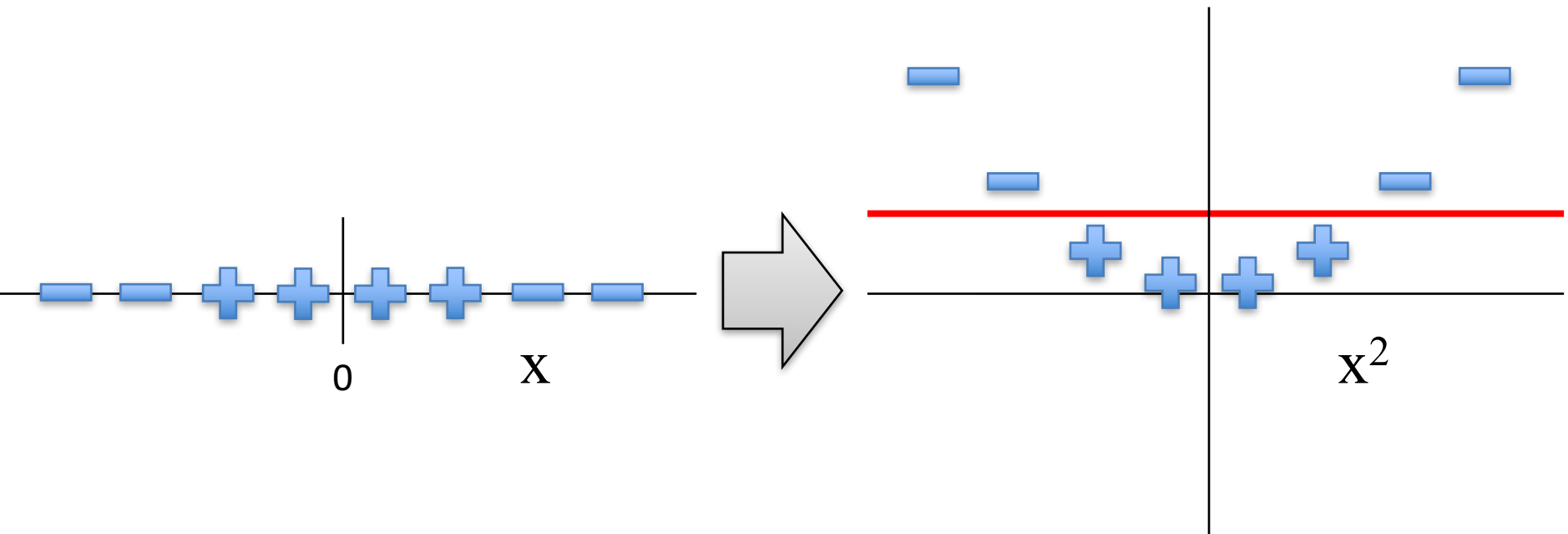
Retinal projection

- Noise in the observation x
  - Measurement error
  - Instrument limitations

# Why Might Predictions be Wrong?

- True non-determinism

- Partial observability
  - hard, soft

- Representational bias

- Algorithmic bias

- Bounded resources

# Representational Bias

- Having the right features ($x$) is crucial

# Support Vector Machines

Doing **_Really_** Well with Linear Decision Surfaces

# Strengths of SVMs

- Good generalization
  - in theory
  - in practice
- Works well with few training instances
- Find globally best model
- Efficient algorithms
- Amenable to the kernel trick

# Minor Notation Change

To better match notation used in SVMs

...and to make matrix formulas simpler

We will drop using superscripts for the $i^{\text{th}}$ instance

| | | | |
|---|---|---|---|
| $i^{\text{th}}$ instance | $\boldsymbol{x}^{(i)}$ $\Longrightarrow$ | $\mathbf{x}_i$ | **Bold** denotes vector |
| $i^{\text{th}}$ instance label | $y^{(i)}$ $\Longrightarrow$ | $y_i$ | **Non-bold** denotes scalar |
| $j^{\text{th}}$ feature of $i^{\text{th}}$ instance | $x_j^{(i)}$ $\Longrightarrow$ | $x_{ij}$ | |

# Linear Separators

- Training instances

$$\mathbf{x} \in \mathbb{R}^{d+1}, x_0 = 1$$

$$y \in \{-1, 1\}$$

- Model parameters

$$\boldsymbol{\theta} \in \mathbb{R}^{d+1}$$

- Hyperplane

$$\boldsymbol{\theta}^\mathsf{T} \mathbf{x} = \langle \boldsymbol{\theta}, \mathbf{x} \rangle = 0$$
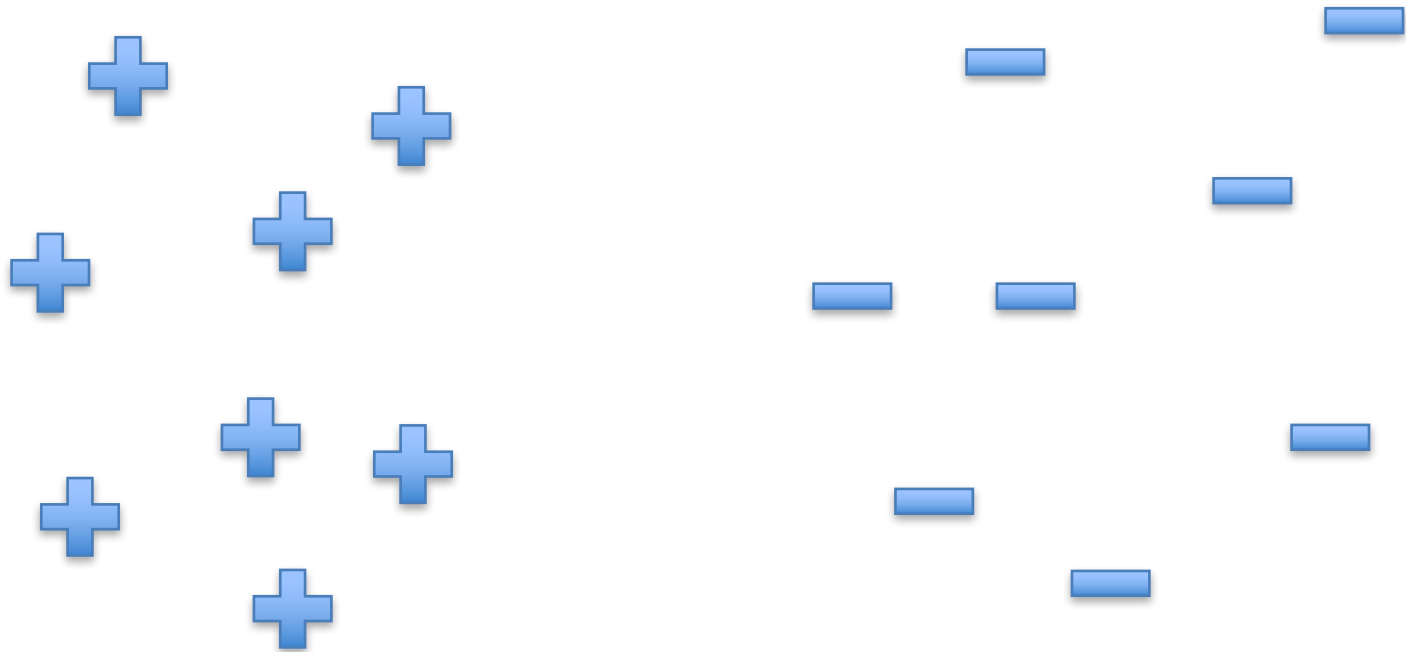
- Decision function

$$h(\mathbf{x}) = \mathrm{sign}(\boldsymbol{\theta}^\mathsf{T} \mathbf{x}) = \mathrm{sign}(\langle \boldsymbol{\theta}, \mathbf{x} \rangle)$$
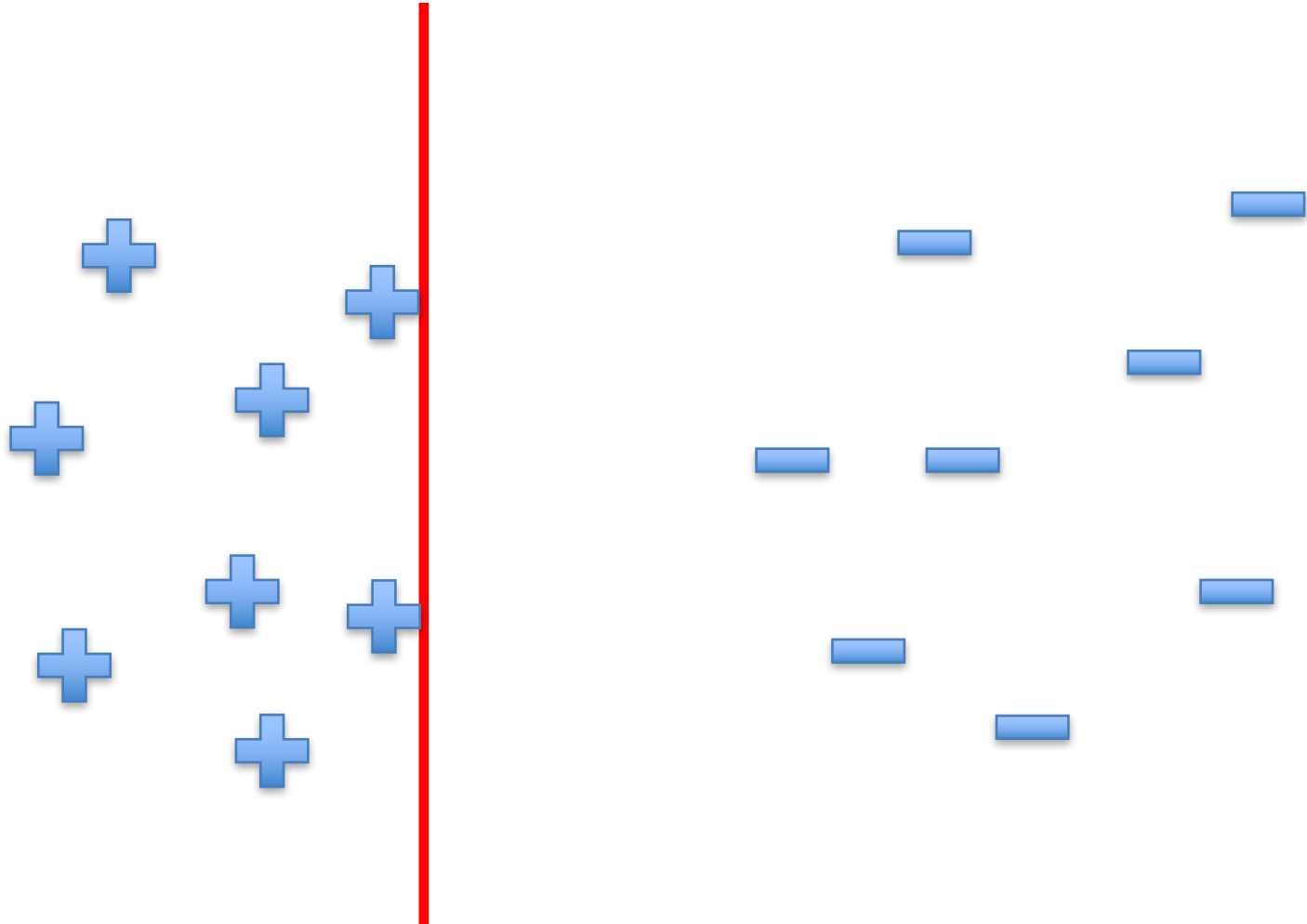
Recall:
Inner (dot) product:

$$\langle \mathbf{u}, \mathbf{v} \rangle = \mathbf{u} \cdot \mathbf{v} = \mathbf{u}^\mathsf{T} \mathbf{v}$$
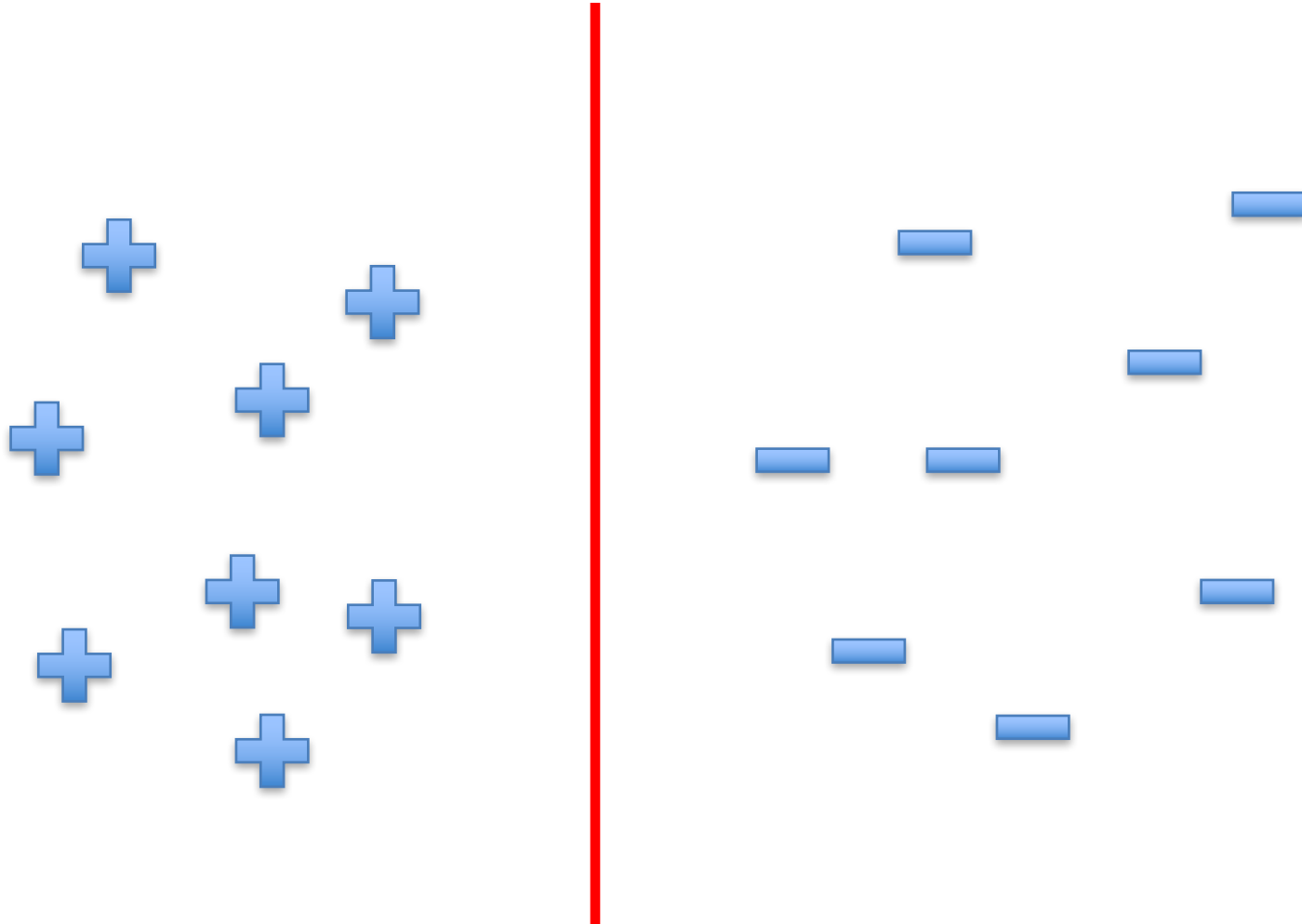
$$= \sum_i u_i v_i$$

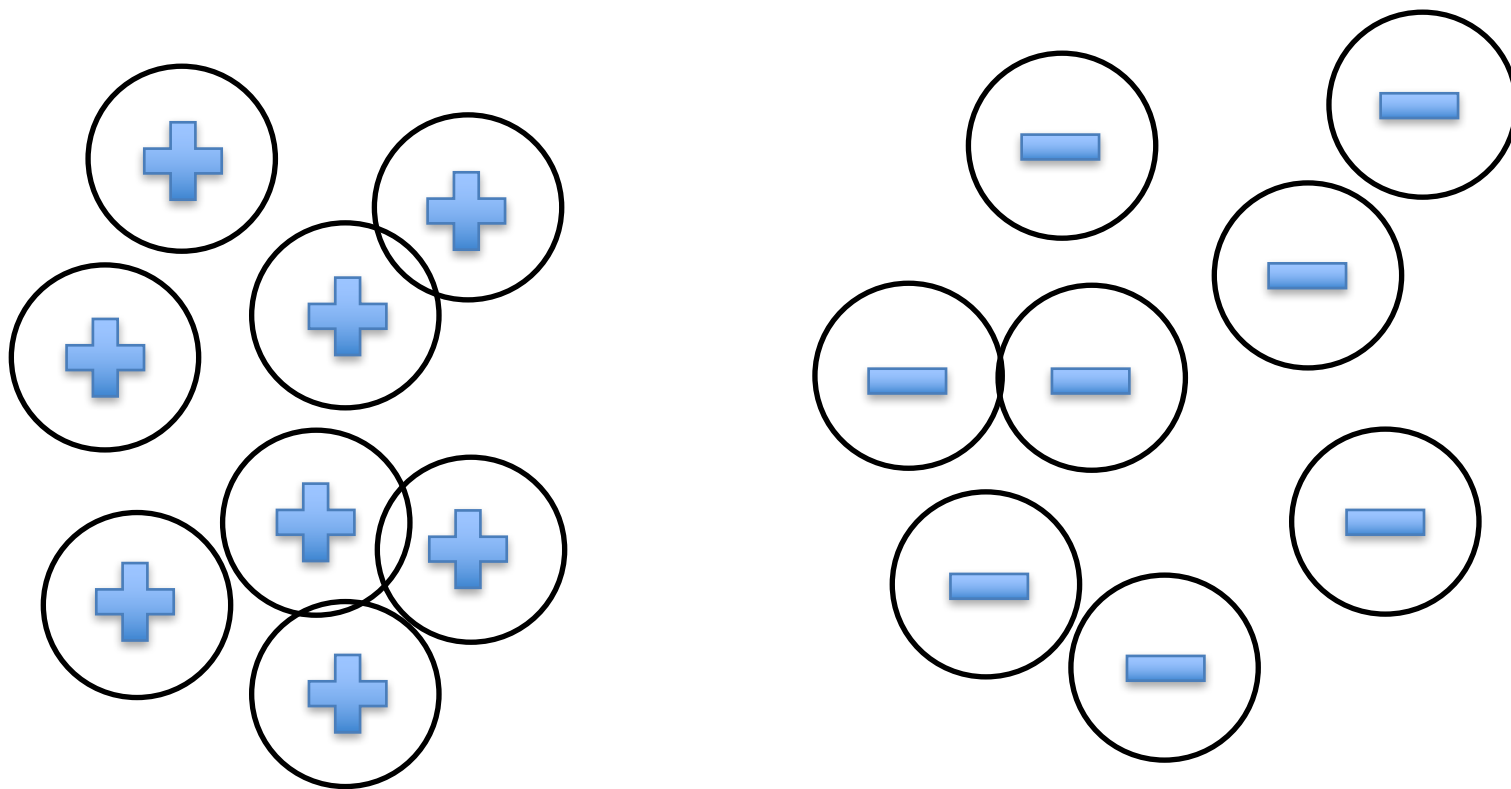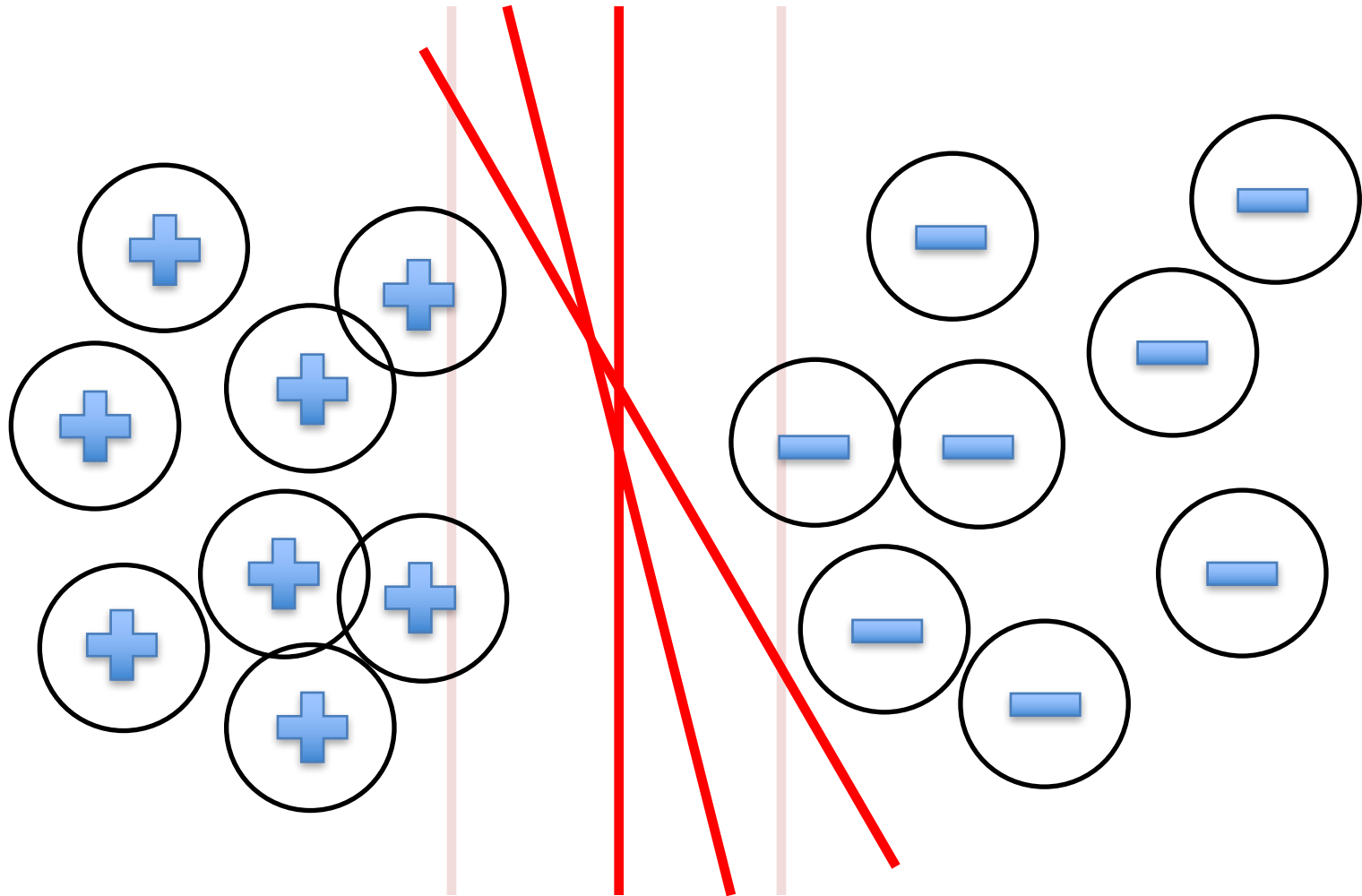# Intuitions

# Intuitions

# Intuitions
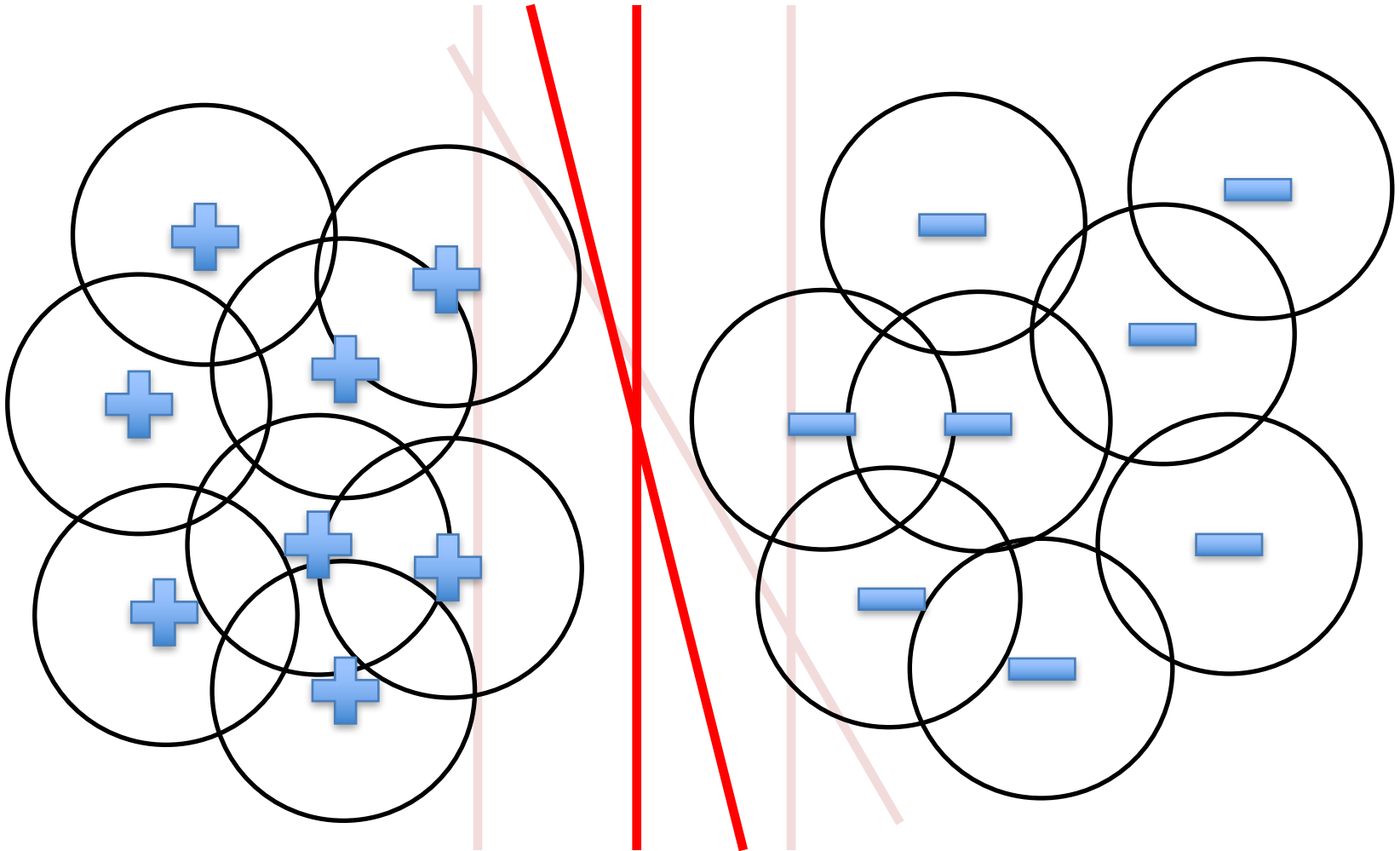
# Intuitions

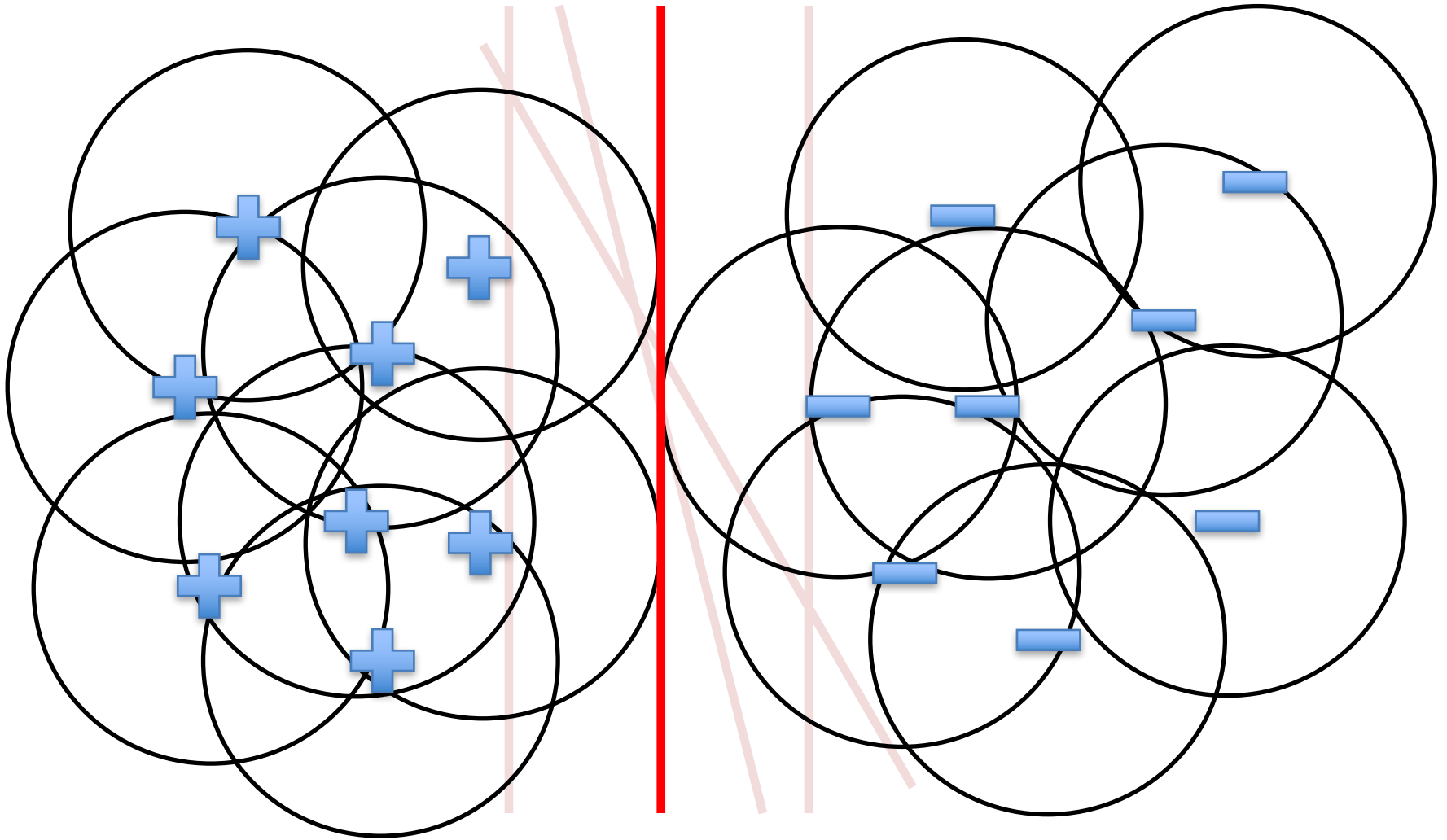# A "Good" Separator
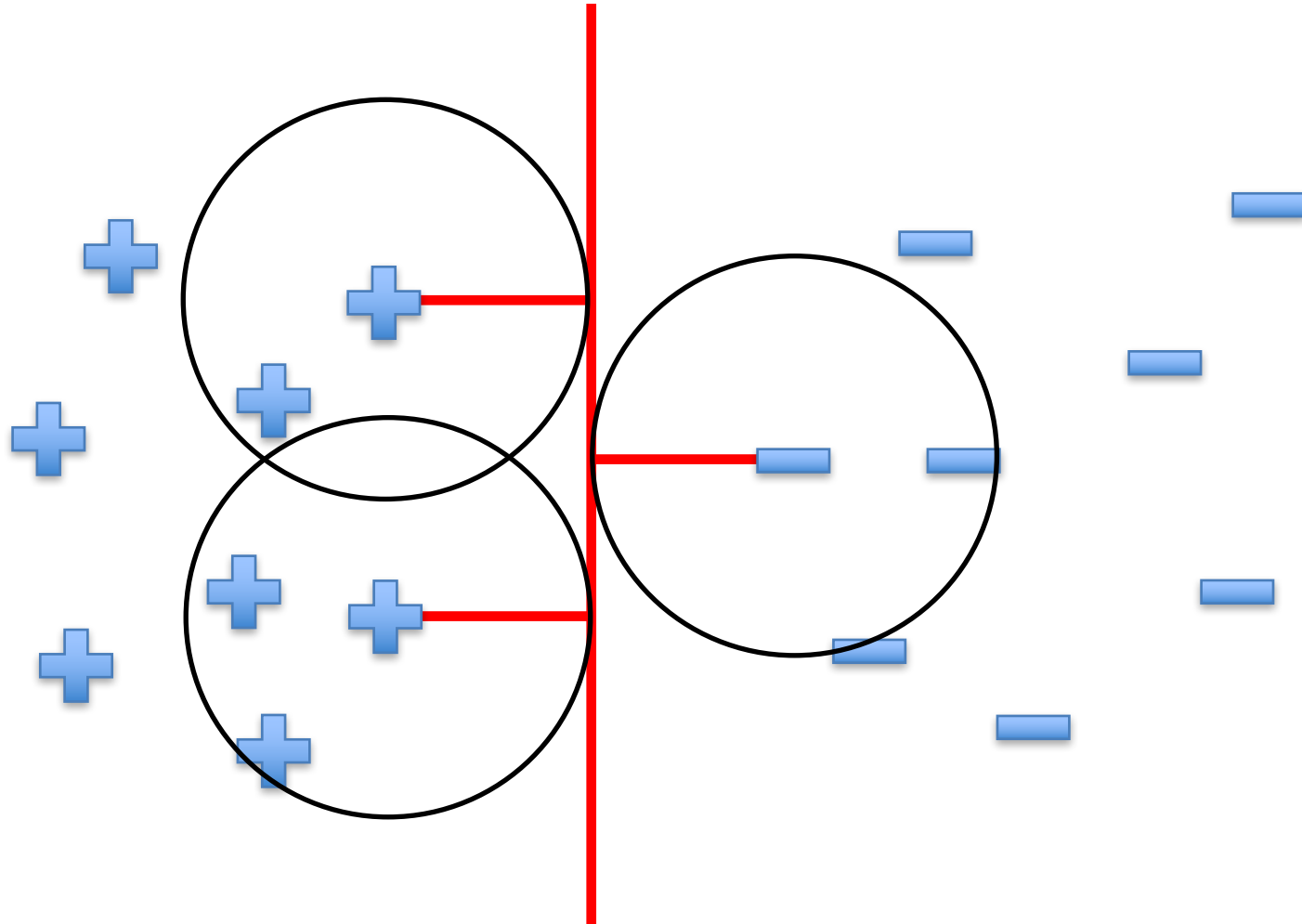
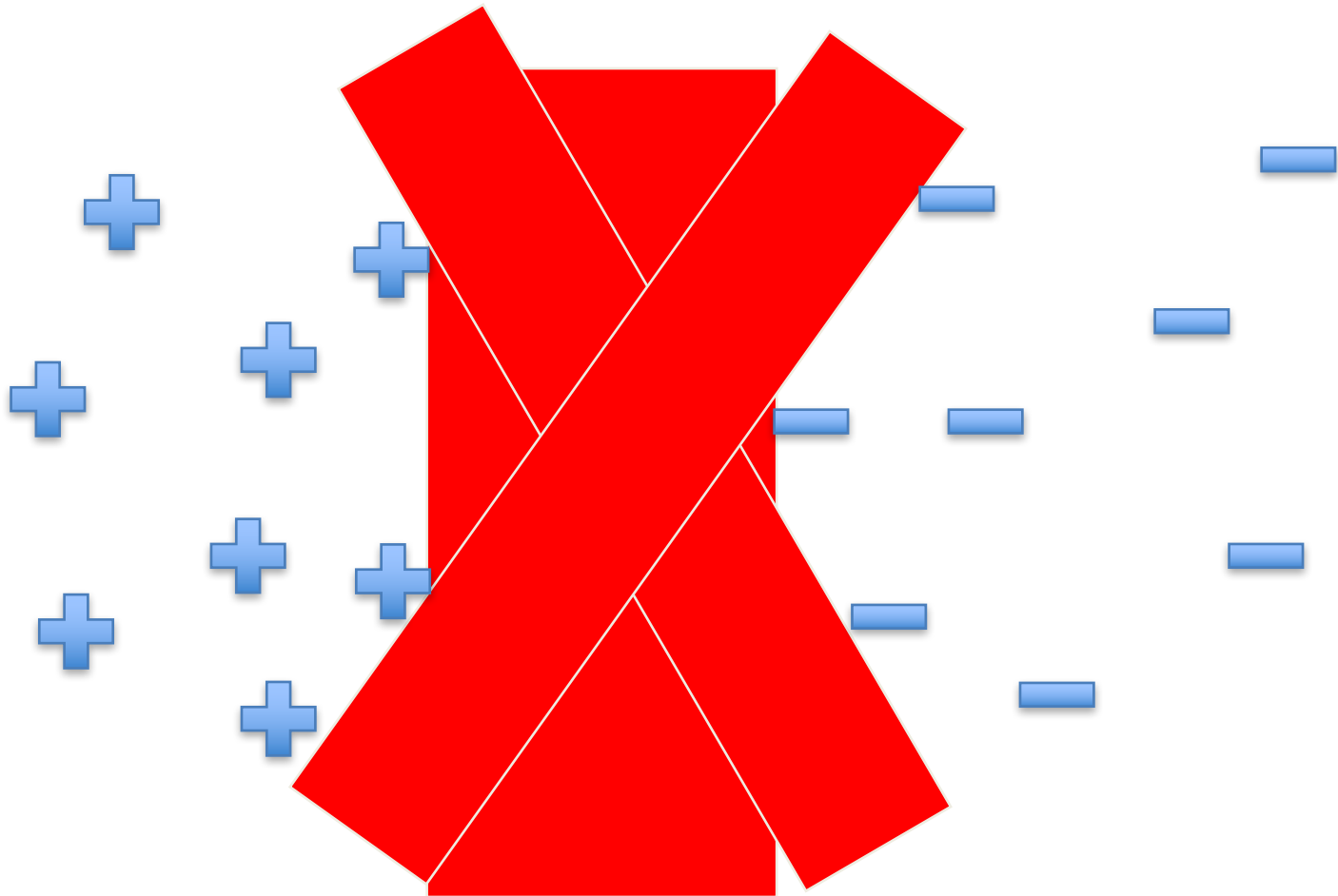# Noise in the Observations

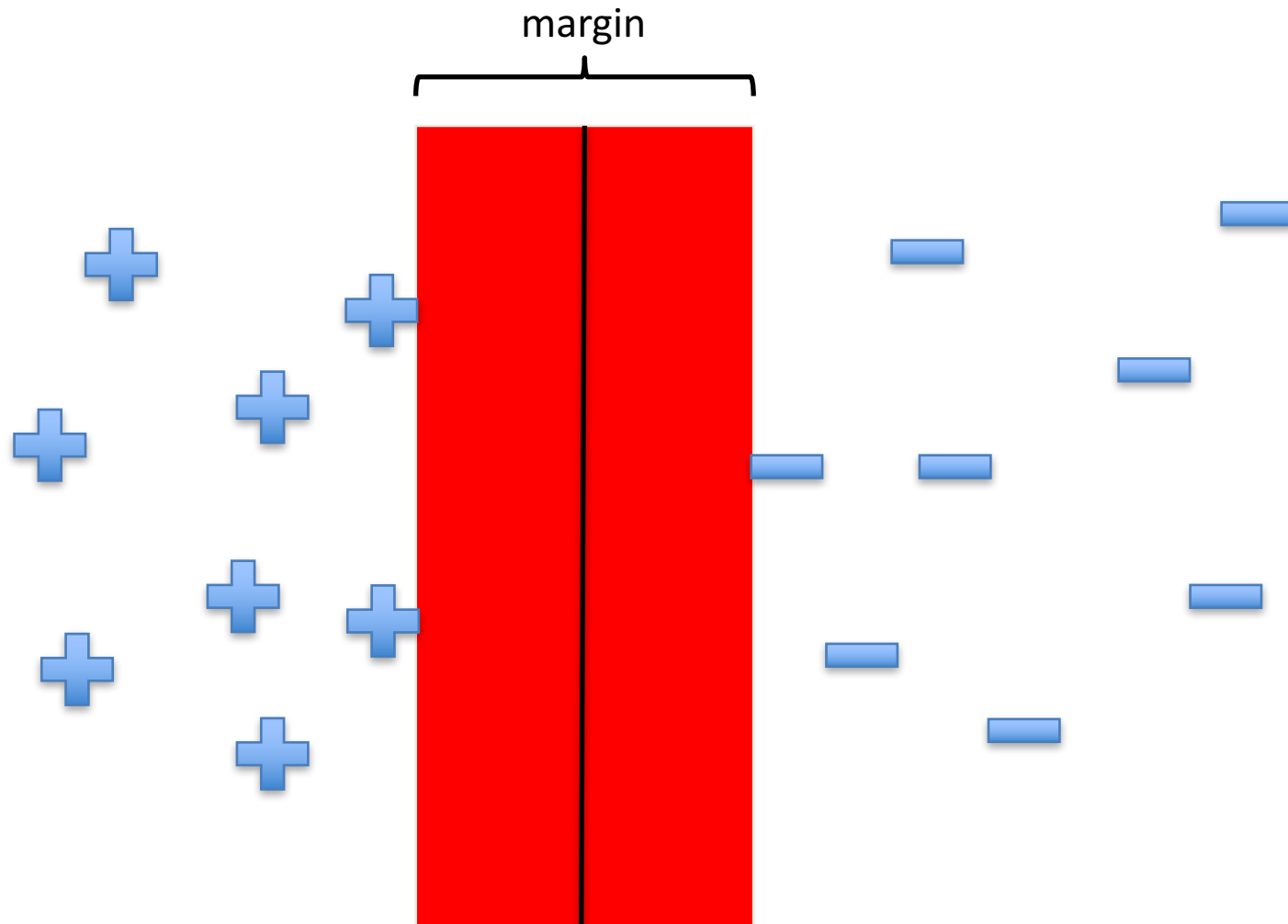# Ruling Out Some Separators

# Lots of Noise

# Only One Separator Remains

# Maximizing the Margin

# "Fat" Separators

# "Fat" Separators


margin

# Why Maximize Margin

Increasing margin reduces *capacity*

- i.e., fewer possible models

**Remember** Lesson from Learning Theory:

- If the following holds:
  - $H$ is sufficiently constrained in size
  - and/or the size of the training data set $n$ is large,

  then low training error is likely to be evidence of low generalization error