

$$\min_{w, b} \frac{1}{2} \|w\|_2^2$$

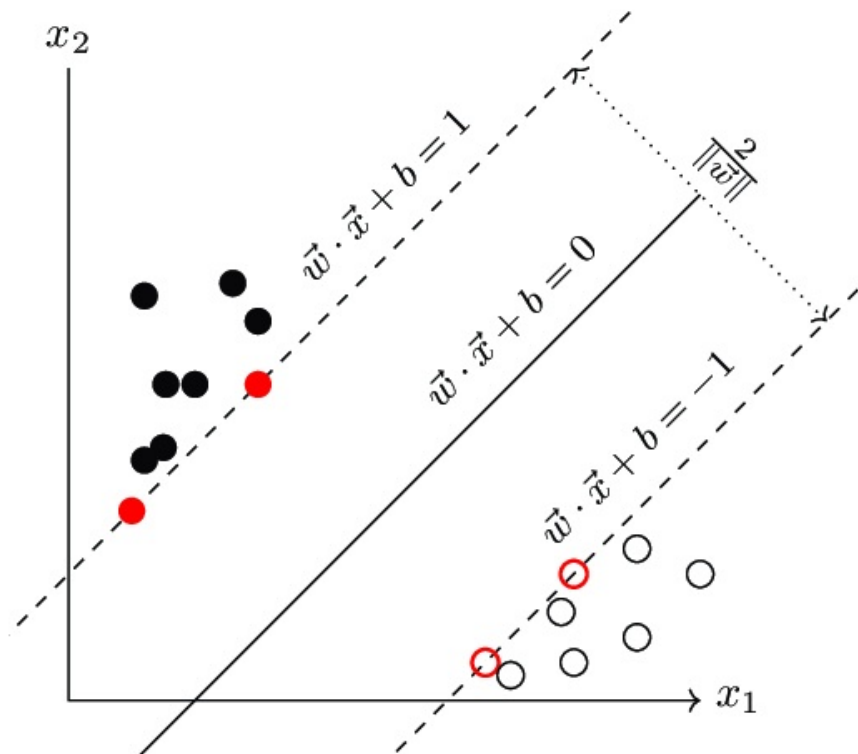
s.t.

$$y_i (\omega^\top x_i + b) \geq 1 \quad i=1, \dots, n$$

$$L(w, b, \lambda)$$

$$\nabla_w L = 0 \Rightarrow \boxed{\omega = \sum_{i=1}^n \lambda_i y_i x_i}$$

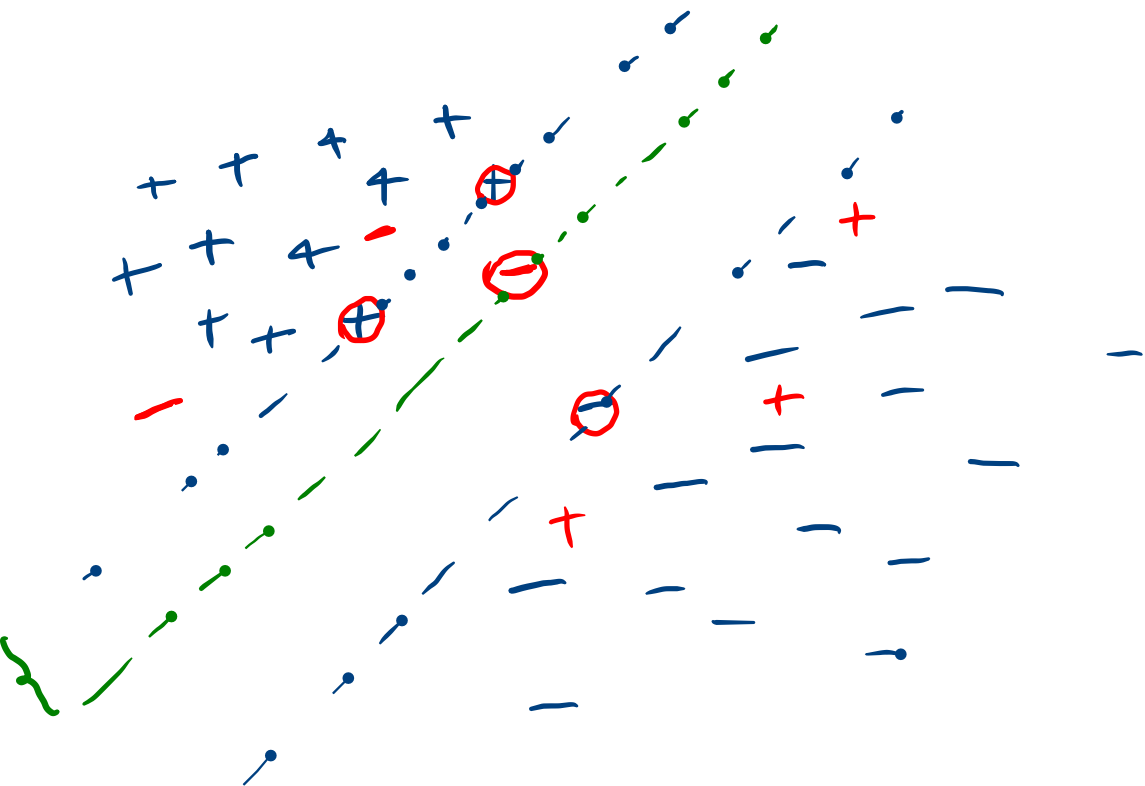
$$\nabla_b L = 0 \Rightarrow \sum_{i=1}^n \lambda_i y_i = 0$$

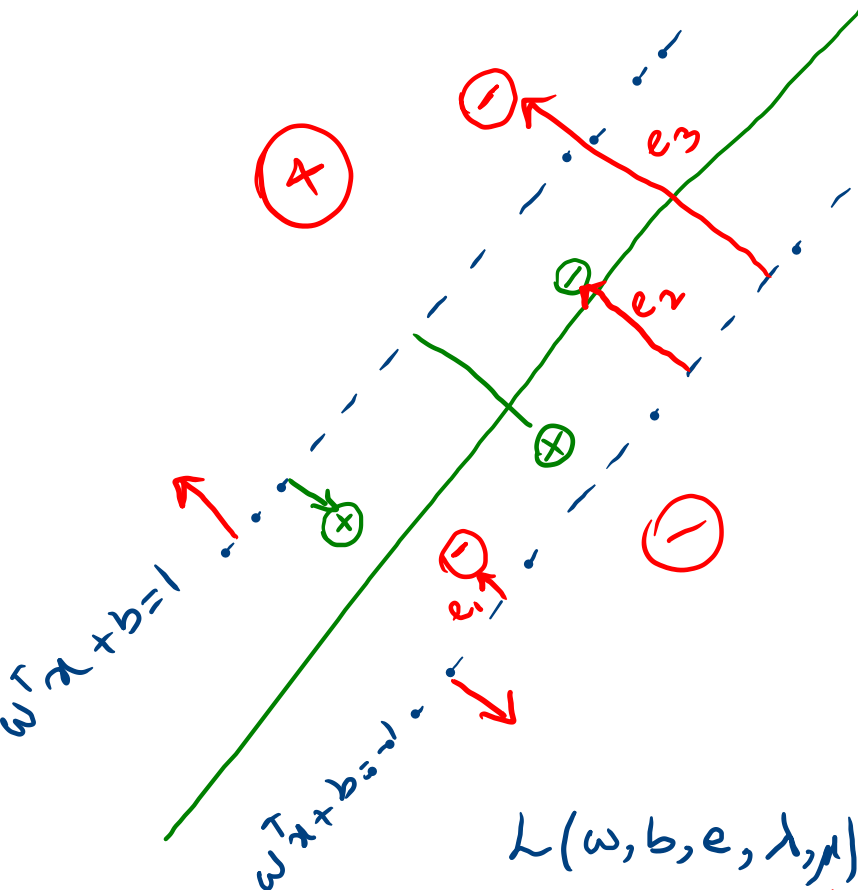


$$\max_{\lambda} \quad -\frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j x_i^T x_j + \sum_{i=1}^n \lambda_i$$

$$\lambda_i \geq 0, \quad \sum \lambda_i y_i = 0$$

Quadratic  
Programming





## Soft Margin SVM

$$\min_{w, b, e} \frac{1}{2} \|w\|_2^2 + c \sum_{i=1}^n e_i$$

s.t.

$$y_i (w^T x_i + b) \geq 1 - e_i$$

$$e_i \geq 0$$

$$L(w, b, e, \lambda, \mu) = \frac{1}{2} \|w\|_2^2 + c \sum_{i=1}^n e_i + \sum_{i=1}^n \lambda_i (1 - e_i - y_i (w^T x_i + b)) + \sum_{i=1}^n \mu_i (-e_i)$$

dual  
var

$$\nabla_{e_j} L = c - \lambda_j - \mu_j = 0 \Rightarrow \boxed{\mu_j = c - \lambda_j} \geq 0 \Rightarrow 0 \leq \lambda_j \leq c$$

$$\max_{\lambda} \quad -\frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j \underbrace{x_i^T x_j}_{\text{red underline}} + \sum_i \lambda_i$$

$$\sum_{i=1}^n \lambda_i y_i = 0, \quad 0 \leq \underbrace{\lambda_i}_{\text{red underline}} \leq C$$

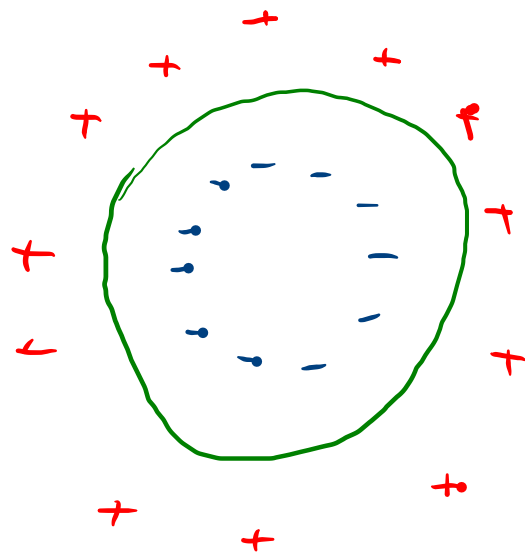
$\lambda^*$

$$K(x_i, x_j) = \Phi(x_i)^T \Phi(x_j)$$

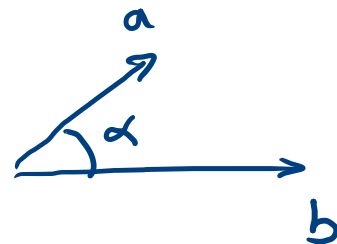
$$\omega = \sum_{i=1}^n \lambda_i y_i \Phi(x_i)$$

$$\omega^T \Phi(x) + b$$

$$g(x) = \sum_{i=1}^n \lambda_i y_i \underbrace{\Phi(x_i)^T \Phi(x)}_{K(x_i, x)} + b \quad \begin{matrix} + \\ \sum \\ - \end{matrix} 0$$



$r$   
 $\theta$



$$|a| |b| \underbrace{\cos \alpha}_1$$

$\Phi$

$x_1, \dots, x_n$

$\underbrace{\Phi(x_1), \dots, \Phi(x_n)}$

Kernel Trick

$$k(x_i, x_j) = \underbrace{\Phi(x_i)^T \Phi(x_j)}$$



$$\cancel{k(x_i, x_j) = x_i x_j + \cos(x_i + x_j)}$$

Mercer's Theorem:

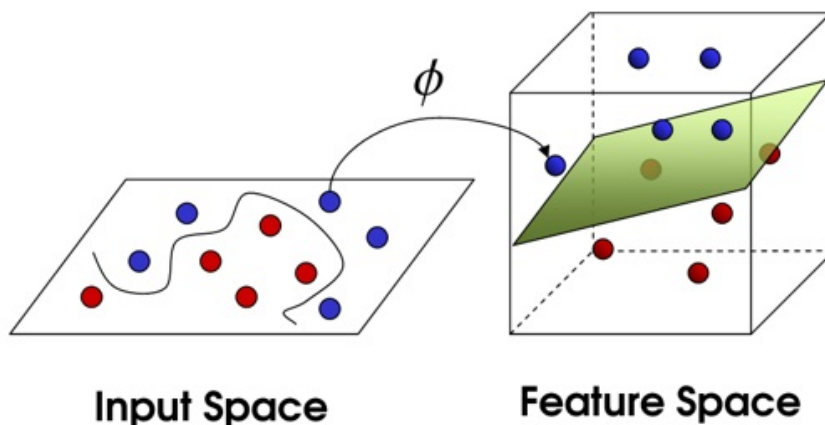
SS

.

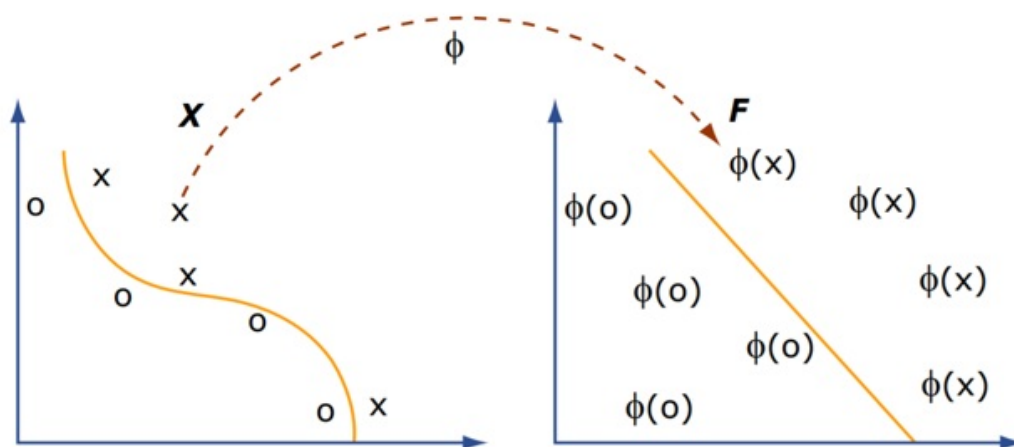




# Transforming the Data



- Feature space is of **higher dimension** than the **input space** in practice
- Computation in the feature space can be **costly** because it is high dimensional



- The **kernel trick** comes to rescue



# kernel trick

- Suppose  $\phi(\cdot)$  is given as follows

$$\phi\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1x_2)$$

- An **inner product** in the feature space is

$$\langle \phi\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right), \phi\left(\begin{bmatrix} y_1 \\ y_2 \end{bmatrix}\right) \rangle = (1 + x_1y_1 + x_2y_2)^2$$

- So, if we define the kernel function as follows, there is no need to **carry out  $\phi(\cdot)$  explicitly**

$$K(\mathbf{x}, \mathbf{y}) = (1 + x_1y_1 + x_2y_2)^2$$

- This use of kernel function to **avoid carrying out  $\phi(\cdot)$  explicitly** is known as the kernel trick



# Examples of Kernel Functions

- Polynomial kernel with **degree**  $d$

$$K(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^T \mathbf{y} + 1)^d$$

- **Gaussian kernel** with width  $\sigma$

$$\rightarrow K(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|^2 / (2\sigma^2))$$

RBF  $\Phi$

- Closely related to radial basis function neural networks
- The feature space is **infinite-dimensional** (it still be written as a dot product in a new feature space  $k(\mathbf{x}, \mathbf{x}_0) = \Phi(\mathbf{x}) * \Phi(\mathbf{x}_0)$ , only with an **infinite number of dimensions**)

- Sigmoid with parameter  $\kappa$  and  $\theta$

$$K(\mathbf{x}, \mathbf{y}) = \tanh(\underline{\kappa} \mathbf{x}^T \mathbf{y} + \underline{\theta})$$

- It does not satisfy the Mercer condition on all  $\kappa$  and  $\theta$

$$k(x, y) = e^{-\|x-y\|^2}$$

$$= e^{-(x-y)^T(x-y)} = \underbrace{e^{-x^T x}}_{\text{}} \underbrace{e^{-y^T y}}_{\text{}} \underbrace{e^{+2x^T y}}_{\text{}}$$

$$e^x = 1 + x + \frac{x^2}{2!} + \dots$$

$$\phi(x)^T \phi(y)$$

## Mercer Theorem

$$\iint \underbrace{g(x)}_{\downarrow} \underbrace{K(x, y)}_{\downarrow} \underbrace{g(y)}_{\downarrow} dx dy \geq 0$$

## Techniques for Constructing New Kernels.

Given valid kernels  $k_1(\mathbf{x}, \mathbf{x}')$  and  $k_2(\mathbf{x}, \mathbf{x}')$ , the following new kernels will also be valid:

$c > 0$

$$\longrightarrow k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}') \quad (6.13)$$

$$k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x}) \underline{k_1(\mathbf{x}, \mathbf{x}')} f(\mathbf{x}') \quad (6.14)$$

$$k(\mathbf{x}, \mathbf{x}') = q(k_1(\mathbf{x}, \mathbf{x}')) \quad (6.15)$$

$$k(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}')) \quad (6.16)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}') \quad (6.17)$$

$$\longrightarrow k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}') \quad (6.18)$$

$$k(\mathbf{x}, \mathbf{x}') = k_3(\phi(\mathbf{x}), \phi(\mathbf{x}')) \quad (6.19)$$

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{A} \mathbf{x}' \quad (6.20)$$

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a) + k_b(\mathbf{x}_b, \mathbf{x}'_b) \quad (6.21)$$

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a)k_b(\mathbf{x}_b, \mathbf{x}'_b) \quad (6.22)$$

where  $c > 0$  is a constant,  $f(\cdot)$  is any function,  $q(\cdot)$  is a polynomial with nonnegative coefficients,  $\phi(\mathbf{x})$  is a function from  $\mathbf{x}$  to  $\mathbb{R}^M$ ,  $k_3(\cdot, \cdot)$  is a valid kernel in  $\mathbb{R}^M$ ,  $\mathbf{A}$  is a symmetric positive semidefinite matrix,  $\mathbf{x}_a$  and  $\mathbf{x}_b$  are variables (not necessarily disjoint) with  $\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)$ , and  $k_a$  and  $k_b$  are valid kernel functions over their respective spaces.

Chapter 6

Bishop

$k_1$

$k_2$

$k(\mathbf{x}, \mathbf{y})$

$k(\phi(\mathbf{x}), \phi(\mathbf{y}))$

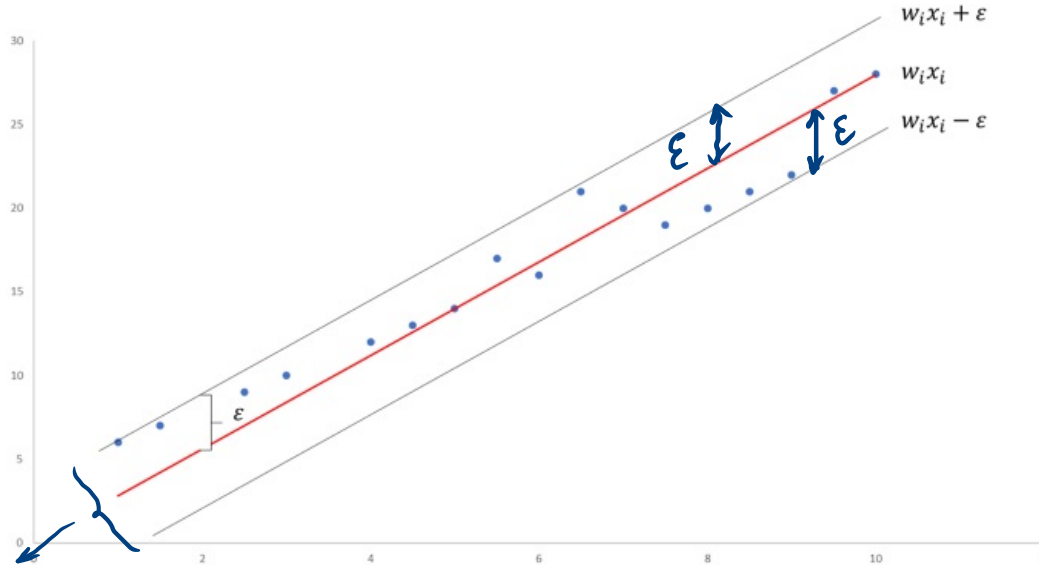
$$K(x, y) = x^T A y = x^T B^T B y = \underbrace{(Bx)^T}_{\phi(x)^T} \underbrace{(By)}_{\phi(y)}$$

$$A = B^T B$$

-



# Support Vector Regression (SVR)



Minimize:

$$\text{MIN } \frac{1}{2} \|\underline{w}\|^2$$

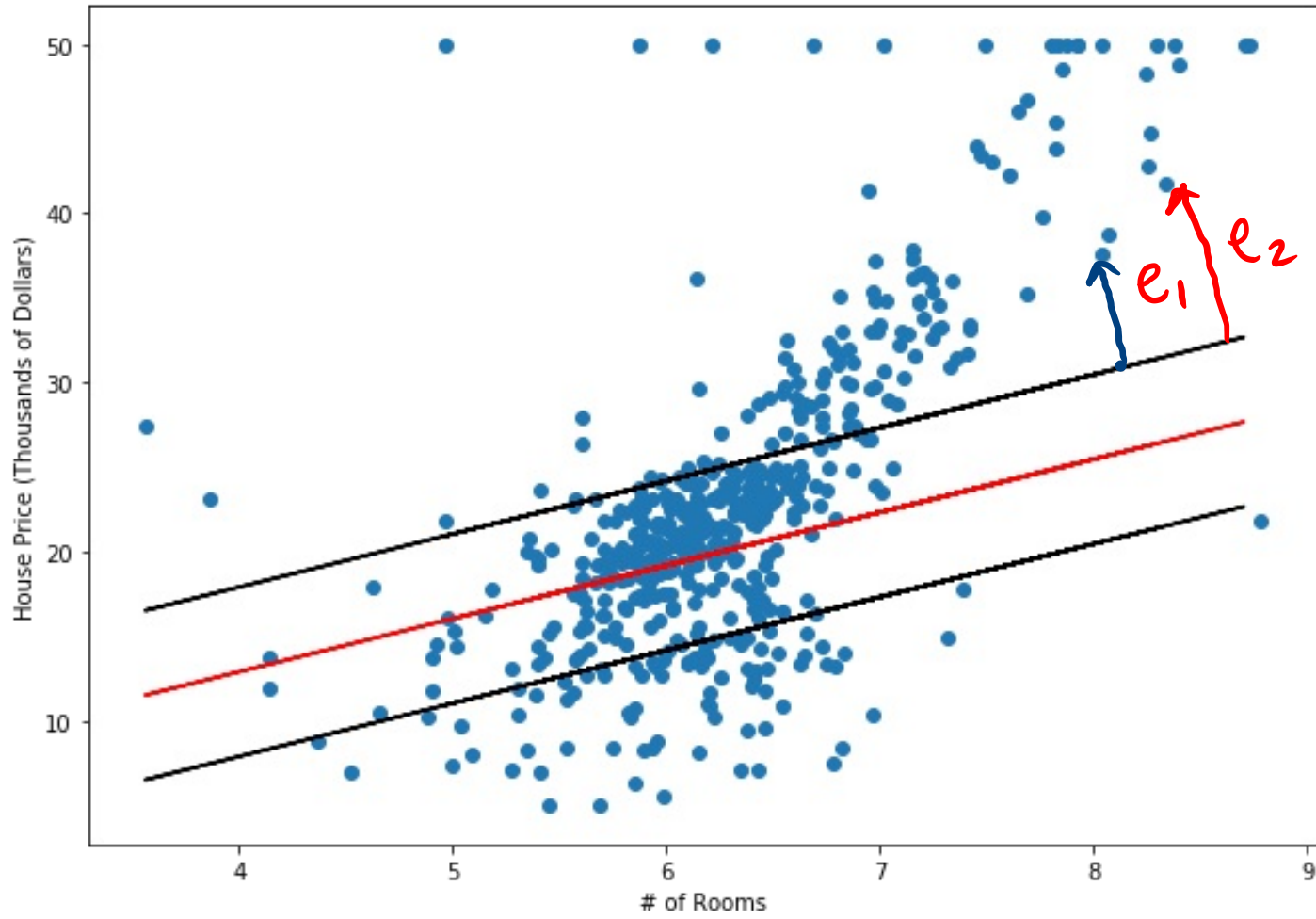
Constraints:

$$|y_i - w_i x_i| \leq \epsilon$$

Ordinary Least Squares (OLS):

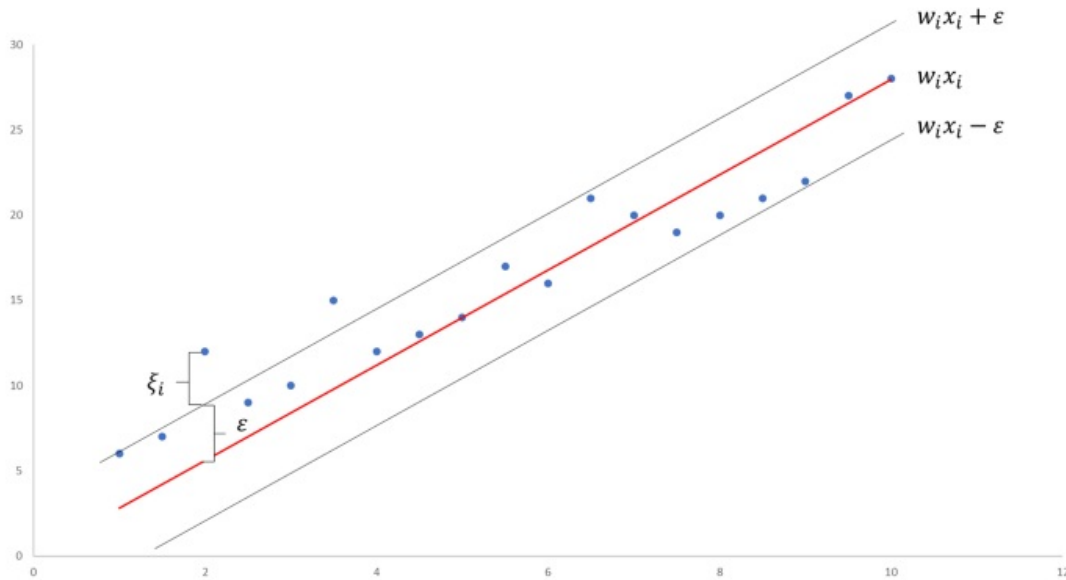
$$\text{MIN } \sum_{i=1}^n (y_i - w_i x_i)^2$$

SVR Prediction



$e_i$

## Using Slack Variables



Minimize:

$$\text{MIN } \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n |\xi_i|$$

Constraints:

$$|y_i - w_i x_i| \leq \varepsilon + \underbrace{|\xi_i|}_{e_i}$$