- Kernel functions,
Support vector machines

- Neural Networks , Decision trees. ... non-convex. $\Rightarrow$ local minimum.

-                     .

# KNN on a grid

# Recall Linear Regression

**Setup**

$$x \quad \xrightarrow{\mathbb{R}^D} \quad \boxed{\phi} \quad \xrightarrow{\phi(x)} \quad \boxed{h} \quad \xrightarrow{\hat{y}} \quad \mathbb{R}$$

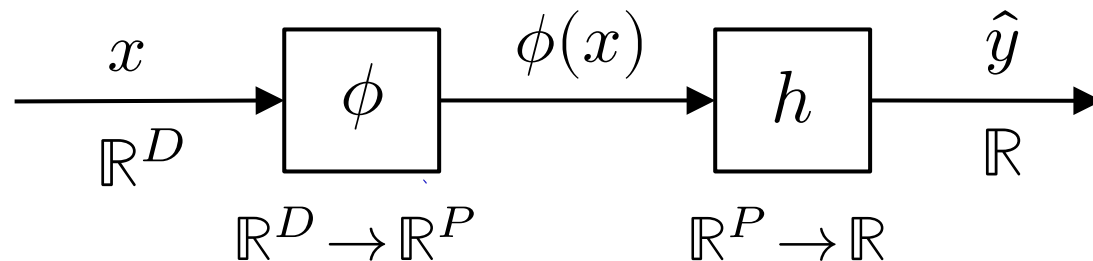$$\mathbb{R}^D \to \mathbb{R}^P \qquad \mathbb{R}^P \to \mathbb{R}$$

Features: $\phi(x) = \begin{bmatrix} \phi_1(x) & \dots & \phi_P(x) \end{bmatrix}$

$x \qquad x^2 \qquad x^3$ (annotations above $\phi_1, \dots, \phi_P$)

$\phi(x)$ ... row vector

$\underline{\theta}_1$ ... column vector

Model: $\boxed{\hat{y} = \theta_0 + \phi(x)\underline{\theta}_1}$

$\phi_1(x) \cdot \theta_1 + \phi_2(x) \cdot \theta_2 \dots + \phi_P(x) \cdot \theta_P$

Goal: $(\hat{\theta}_0, \underline{\hat{\theta}}_1) = \underset{(\theta_0, \underline{\theta}_1) \in \mathbb{R}^{P+1}}{\operatorname{argmin}} \sum_{i=1}^{N} \left( \theta_0 + \phi(x_i)\,\underline{\theta}_1 - y_i \right)^2 + \lambda \sum_{j=1}^{P} \theta_j^2$

Ridge regul.

# Training data



$N \gg P$

| | | | P | | | $\phi(x_i)$ |
|---|---|---|---|---|---|---|
| 0 | 0.247746 | 36.0 | 266. | 0.247746 | 88.019654 | 57.000823 |
| 1 | 0.179340 | 37.0 | 365.0 | 0.179340 | 27.211935 | 22.471227 |
| 2 | 0.956807 | 21.0 | 151.0 | 0.956807 | 97.456012 | 56.357366 |
| 3 | 0.869653 | 43.0 | 437. | 0.869653 | 43.203221 | 34.75 |
| 4 | 0.825345 | 47.0 | 160.0 | 0.825345 | 98.933930 | 64.426163 |
| 5 | 0.331114 | 321.0 | 371.0 | 0.331114 | 8.257917 | 14.218720 |
| 6 | 0.765523 | 17.0 | 364.0 | 0.765523 | 96.696783 | 59.123769 |
| 7 | 0.956807 | 21.0 | 151.0 | 0.956807 | 97.456012 | 52.983470 |

$\Phi$  $Y$  $y_i$

$\hat{\mu}_X$  $\hat{\mu}_Y$

Matrices:

$$\Phi = \begin{bmatrix} \phi(x_1) \\ \vdots \\ \phi(x_N) \end{bmatrix} \in \mathbb{R}^{N \times P} \qquad \mathbf{Y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} \in \mathbb{R}^{N \times 1}$$

Means:

$$\hat{\mu}_X = \frac{1}{N} \mathbf{1}_N^T \Phi \qquad\qquad \hat{\mu}_Y = \frac{1}{N} \mathbf{1}_N^T \mathbf{Y}$$
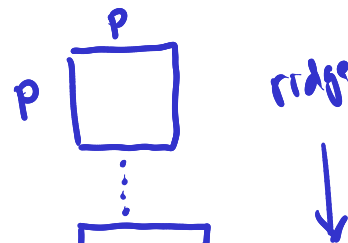
Centered matrices:

$$\Phi_c = \Phi - \mathbf{1}_N \hat{\mu}_X \qquad\qquad \mathbf{Y}_c = \mathbf{Y} - \mathbf{1}_N \hat{\mu}_Y$$

## **Solution**

P



$$P$$

ridge

**Stationarity condition:** $\left(\Phi_c^T \Phi_c + \lambda I_P\right) \widehat{\underline{\theta}}_1 = \Phi_c^T \mathbf{Y}_c$

P x P

**Optimal parameters:** $\cdot \ \widehat{\underline{\theta}}_1 = \left(\Phi_c^T \Phi_c + \lambda I_P\right)^{-1} \Phi_c^T \mathbf{Y}_c$   normal eqn.

$\cdot \ \widehat{\theta}_0 = \hat{\mu}_Y - \hat{\mu}_X \ \widehat{\underline{\theta}}_1$

**Prediction:** $\hat{y} = \widehat{\theta}_0 + \phi(x)\widehat{\underline{\theta}}_1$

$h(x).$

$\qquad = \hat{\mu}_Y + \phi_c(x)\widehat{\underline{\theta}}_1$
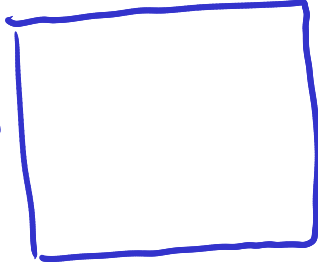
another prediction formula.

# **Solution**

Stationarity condition: $\left(\Phi_c^T\Phi_c + \lambda I_P\right)\hat{\underline{\theta}}_1 = \Phi_c^T\mathbf{Y}_c$

Optimal parameters: $\hat{\underline{\theta}}_1 = \left(\Phi_c^T\Phi_c + \lambda I_P\right)^{-1}\Phi_c^T\mathbf{Y}_c$

$\hat{\theta}_0 = \hat{\mu}_Y - \hat{\mu}_X\,\hat{\underline{\theta}}_1$

Prediction: $\hat{y} = \hat{\theta}_0 + \phi(x)\hat{\underline{\theta}}_1$

$\hat{y} = \hat{\mu}_Y - \hat{\mu}_X\,\hat{\underline{\theta}}_1 + \phi(x)\hat{\underline{\theta}}_1$

$\hat{y} - \hat{\mu}_Y = (\phi(x) - \hat{\mu}_X)\hat{\underline{\theta}}_1$

$\hat{y}_c = \phi_c(x)\hat{\underline{\theta}}_1$

# Kernel form of Ridge regression
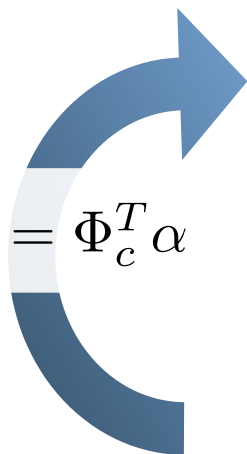
$N$

$\Phi_c \Phi_c^T = N$
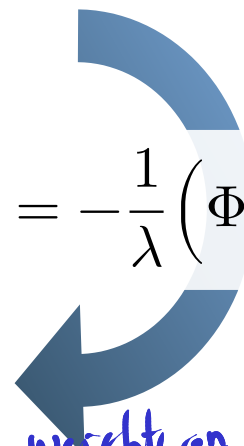
## Standard form

$P \times P$

$$\underline{\hat{\theta}}_1 = \left(\Phi_c^T \Phi_c + \lambda I_P\right)^{-1} \Phi_c^T \mathbf{Y}_c$$

$$\hat{y} = \hat{\mu}_Y + \phi_c(x)\underline{\hat{\theta}}_1$$

$$\underline{\hat{\theta}}_1 = \Phi_c^T \alpha$$

$$\alpha = -\frac{1}{\lambda}\left(\Phi_c\,\underline{\hat{\theta}}_1 - \mathbf{Y}_c\right)$$

## Kernel form

$N \times N$

$$\alpha = \left(\Phi_c \Phi_c^T + \lambda\, I_N\right)^{-1} \mathbf{Y}_c \quad \dots \; \mathbb{R}^N \; \dots \; \text{weights on each training data point.}$$

$$\hat{y} = \hat{\mu}_Y + \sum_{i=1}^{N} \alpha_i \phi_c(x_i) \cdot \phi_c(x)$$

$K(x_i, x)$.

# Kernel form of Ridge regression

**Training:**

$$\alpha = \left( \Phi_c \Phi_c^T + \lambda I_N \right)^{-1} \mathbf{Y}_c$$

N×N

$$\mathbb{K} = \Phi_c \Phi_c^T \quad \in \mathbb{R}^{N \times N} \quad \text{... Kernel matrix}$$

$$= \begin{bmatrix} \phi_c(x_1) \cdot \phi_c(x_1) & \dots & \phi_c(x_1) \cdot \phi_c(x_N) \\ \vdots & & \vdots \\ \phi_c(x_N) \cdot \phi_c(x_1) & \dots & \phi_c(x_N) \cdot \phi_c(x_N) \end{bmatrix}$$

$$= \begin{bmatrix} k(x_1, x_1) & \dots & k(x_1, x_N) \\ \vdots & & \vdots \\ k(x_N, x_1) & \dots & k(x_N, x_N) \end{bmatrix}$$

... Kernel function

$$k(x, z) = \phi_c(x) \cdot \phi_c(z)$$

symmetric

symmetric

$$\begin{bmatrix} \phi_c(x_1) \\ \phi_c(x_2) \\ \vdots \\ \phi_c(x_N) \end{bmatrix} \begin{bmatrix} \phi_c^T(x_1) & \dots & \phi_c^T(x_N) \end{bmatrix}$$

$$\Phi_c^T$$

$$\Phi_c$$

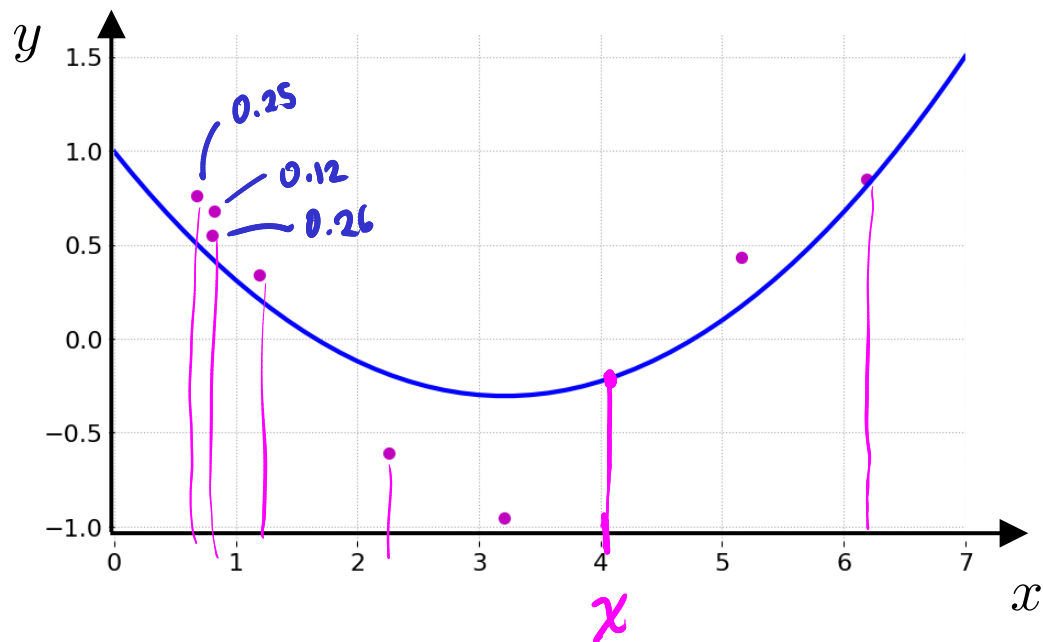**Prediction:** $h(x): \hat{y} = \hat{\mu}_Y + \sum_{i=1}^{N} \alpha_i \, k(x_i, x)$

# Example

$$N = 8$$

$$D = 1$$

$$P = 2$$

$$\phi(x) = (x, x^2)$$

$$\Phi_c = \begin{bmatrix} \overset{x-\text{mean}}{-1.87} & \overset{x^2-\text{mean}}{-9.98} \\ -1.74 & -9.80 \\ -1.72 & -9.77 \\ -1.35 & -9.02 \\ -0.28 & -5.33 \\ 0.8 & -0.13 \\ 2.62 & 16.15 \\ 3.66 & 27.88 \end{bmatrix} \qquad \mathbf{Y}_c = \begin{bmatrix} 0.50 \\ 0.29 \\ 0.42 \\ 0.09 \\ -0.87 \\ -1.21 \\ 0.18 \\ 0.60 \end{bmatrix}$$



$$X = \begin{bmatrix} 0.8 \\ 0.9 \\ 1.1 \\ \vdots \\ 6.2 \end{bmatrix} \xrightarrow{\phi} \begin{bmatrix} \overset{x}{0.8} & \overset{x^2}{0.64} \\ 0.9 & 0.81 \\ 1.1 & 1.21 \\ \vdots & \vdots \\ \underset{\overline{(\overline{x})}}{6.2} & \underset{\overline{(\overline{x})}}{\phantom{x}} \end{bmatrix} \quad \hat{\mu}_x$$

## Standard solution

$$\hat{\theta}_0 = 1$$

$$\underline{\hat{\theta}}_1 = \begin{bmatrix} -0.81 \\ 0.13 \end{bmatrix} \begin{array}{l} x \\ x^2 \end{array}$$

$$\hat{y}(x) = 1 - 0.81x + 0.13x^2$$

## Kernel-based solution

$$\alpha = \begin{bmatrix} 0.25 \\ 0.12 \\ 0.26 \\ 0.13 \\ -0.42 \\ -0.65 \\ 0.26 \\ 0.04 \end{bmatrix} \in \mathbb{R}^8 \quad .4$$

$$\phi(x) = (x, x^2)$$

$$\phi_c(x) = \phi(x) - \hat{\mu}_X = \left( x - \bar{x}, \; x^2 - \bar{x} \right)$$

$$k(x,z) = \phi_c(x) \cdot \phi_c(z)$$

$$\hat{y}(x) = \hat{\mu}_Y + \sum_{i=1}^{N} \alpha_i k(x_i, x) \quad 4$$

# Question

Given an *arbitrary* function $k(x, z)$, are there conditions that guarantee the existence of a feature function $\phi(x)$ such that $k(x, z) = \phi(x) \cdot \phi(z)$?

**Answer:** Yes!

→ 1. $K(x, z)$ must be symmetric: $K(x, z) = K(z, x) \quad \forall x, z \in \mathbb{R}^D$.

→ 2. $K(x, z)$ must be "positive semi-definite".

any training data.

$\mathbb{K} \implies$ must be positive semi-definite.

$x^T \mathbb{K} x \geq 0.$

$\forall x \in \mathbb{R}^D.$

**Example:** Polynomial kernel: $\quad k(x, z) = \left( x^T z + 1 \right)^{\boxed{d}}$

The corresponding feature vector $\phi(x)$ contains monomials of $x$ up to order $d$

$$x = \begin{bmatrix} x^1 \\ \vdots \\ x^D \end{bmatrix}$$

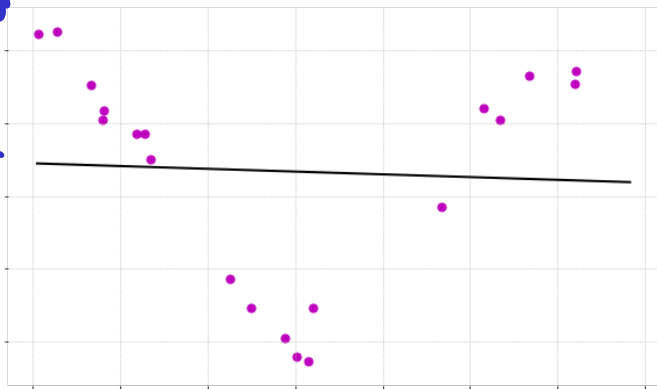$d = 1 : \quad \phi = \left( x^1, x^2, \ldots, x^D \right).$

$d = 2 : \quad \phi = \left( x^1, x^2, \ldots, x^D, (x^1)^2, \ldots, (x^D)^2, x^1 \cdot x^2, x^1 \cdot x^3, \ldots, x^{D-1} x^D \right).$

$d = 3 : \quad \left( \quad , (x^1)^3, \ldots (x^D)^3, x^1 (x^2)^2, \ldots \right).$
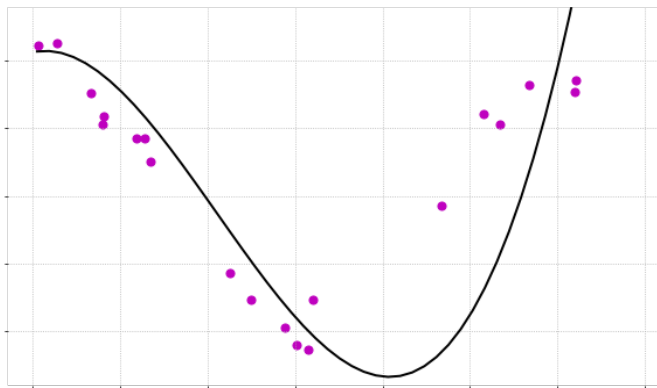
**Linear** $k(x, z) = x^T z$

**Quadratic** $k(x, z) = (x^T z + 1)^2$

**Poly 4** $k(x, z) = (x^T z + 1)^4$
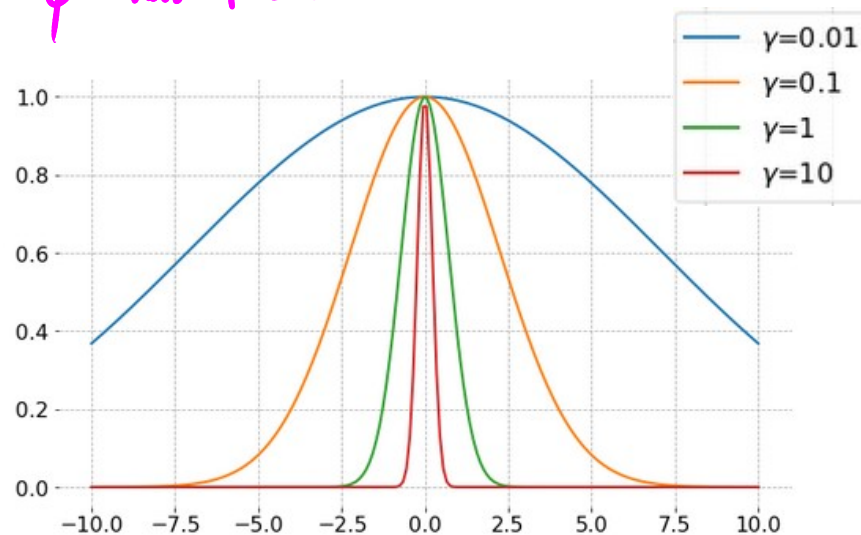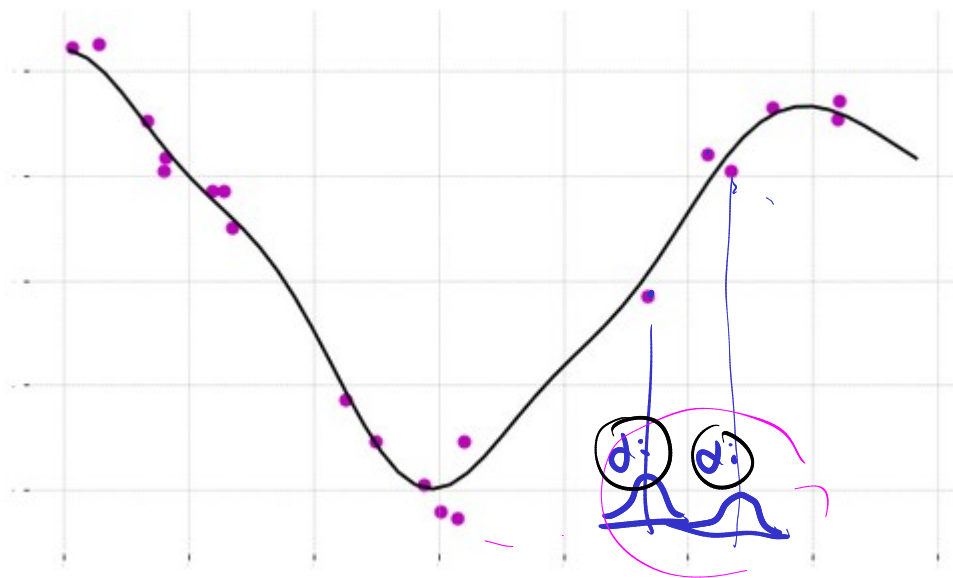
**Poly 10** $k(x, z) = (x^T z + 1)^{10}$

Simple (ridge). linear regression.

······ Overfitting.

# Gaussian kernel – a.k.a. Radial basis function (RBF)

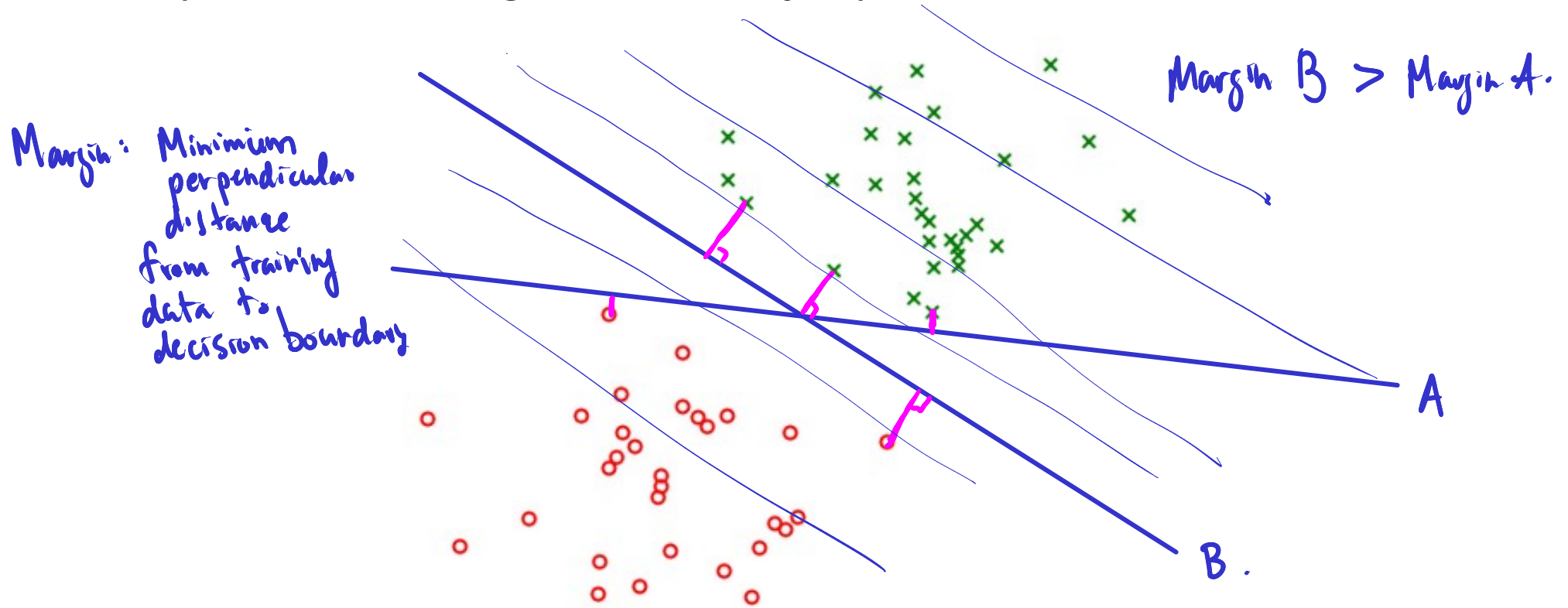$$k(x, z) = \exp\left(-\gamma \left\| x - z \right\|^2\right)$$

$\gamma \approx 1/\sigma^2$

"$\phi$" has $P = \infty$

# Maximum margin classifier

Assumption: The training data is linearly separable

Margin: Minimum perpendicular distance from training data to decision boundary

Margin B > Margin A.

A

B.

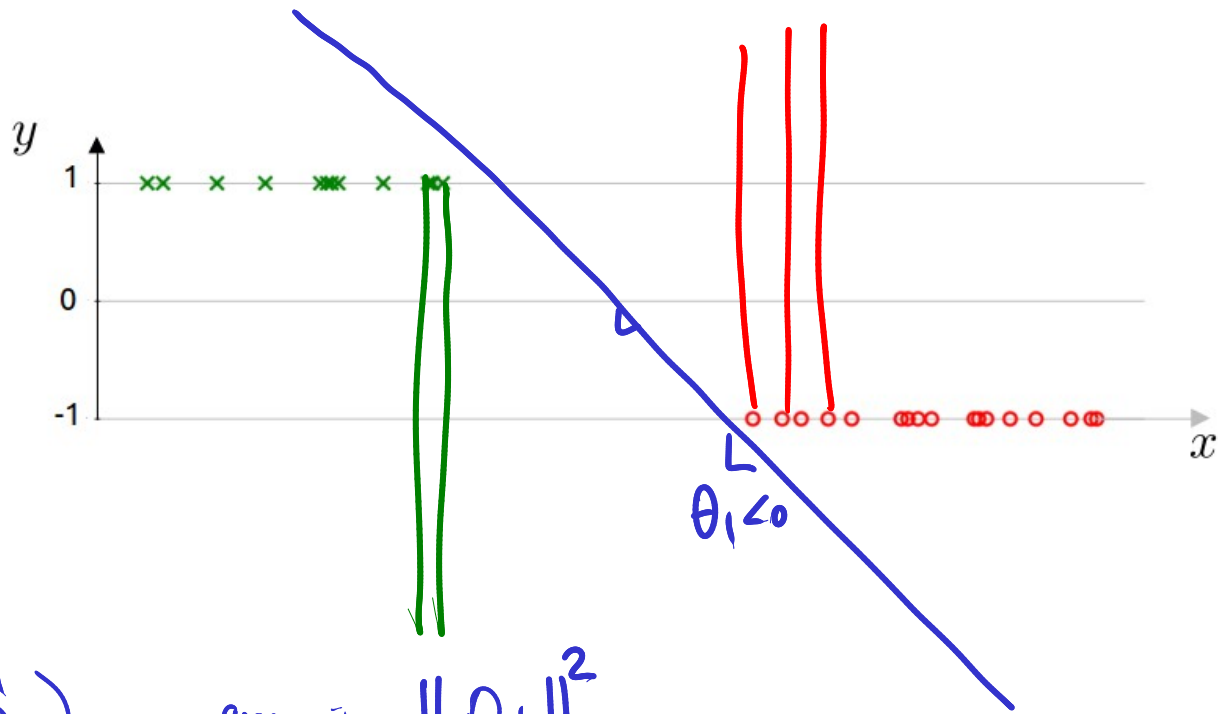Encoding: $\{1, -1\}$.



$\alpha(x) = \theta_0 + \theta_1 x$

$$\left( \hat{\Theta}_0, \hat{\Theta}_1 \right) = \underset{\theta_0, \theta_1}{\text{argmin}} \quad \Theta_1 \qquad \dots \text{convex (linear)}.$$

$$\text{s.t.} \quad y_i = 1 : \quad \alpha(x_i) > 1 \qquad (\text{green})$$

$$y_i = -1 : \quad \alpha(x_i) < -1 \qquad (\text{red})$$

$$\alpha(x_i) = \Theta_0 + \Theta_1 x_i$$

$$\left(\hat{\theta}_0, \hat{\theta}_1\right) = \underset{\theta_0, \theta_1}{\arg\min} \;\left\|\theta_1\right\|^2 \qquad \ldots \text{ convex}$$

$$\text{s.t.} \quad y_i = 1 : \quad \alpha(x_i) \geq 1 \qquad \text{(green)} \qquad \text{(quadratic)}.$$

$$\qquad\qquad y_i = -1 : \quad \alpha(x_i) \leq -1 \qquad \text{(red)}$$
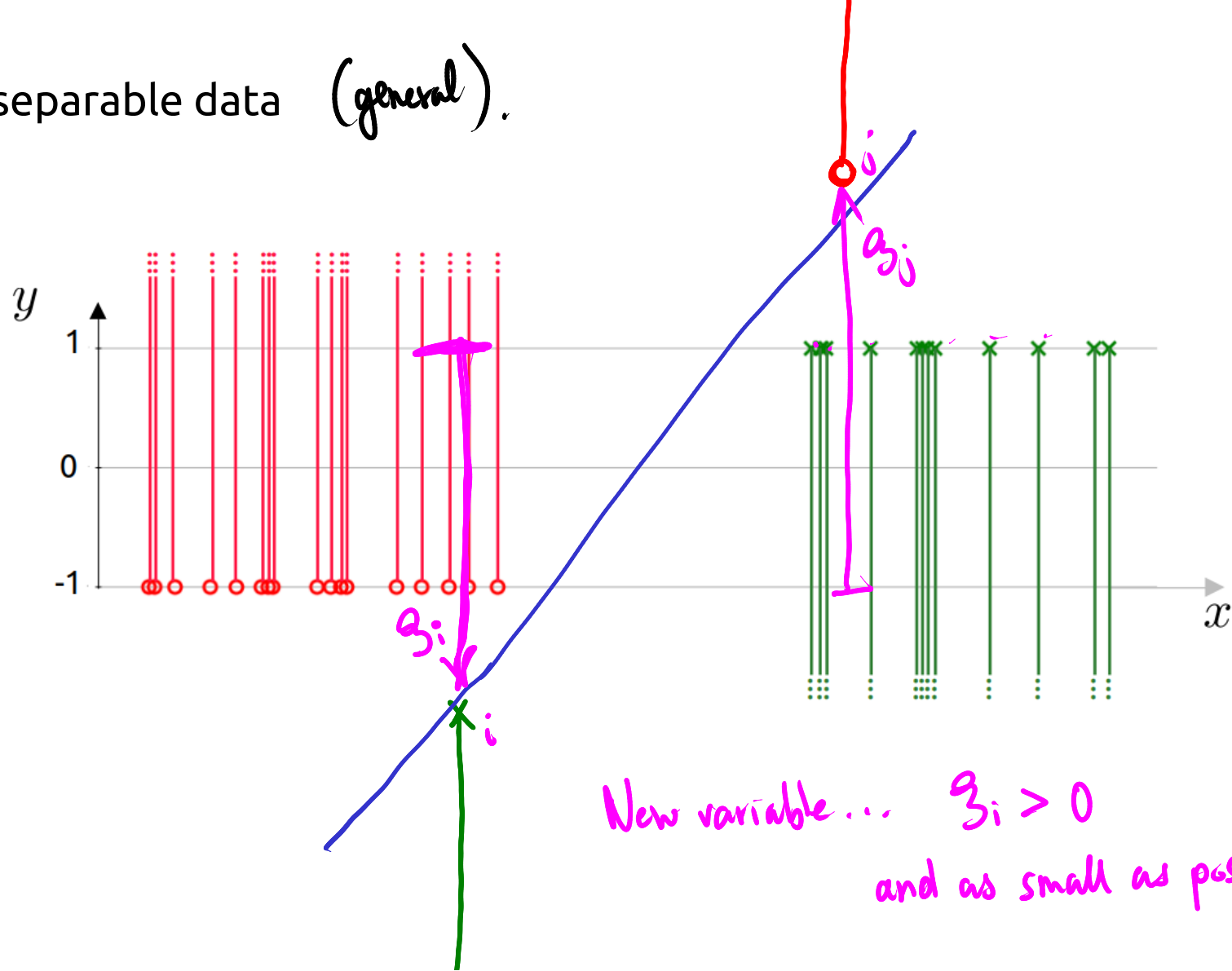
$$\alpha(x_i) = \theta_0 + \theta_1 x_i$$

Simplification: Multiply both sides of the inequality constraints by $y_i$

$$\text{minimize} \quad \| \underline{\theta}_1 \|^2$$

$$\text{s.t.} \quad y_i \, \alpha(x_i) \geq 1$$

Non-separable data (general).



$y$

1

0

-1

$x$

$z_j$

$j$

$z_i$

$x_i$

New variable... $z_i > 0$
and as small as possible.

$$\left( \hat{\theta}_0, \hat{\theta}_1 \right) = \underset{\theta_0, \theta_1, \mathcal{G}_i}{argmin} \left\| \underset{\rightarrow}{\theta_1} \right\|^2 + C \sum_{i=1}^{N} \mathcal{G}_i$$

$$s.t. \qquad \mathcal{G}_i \geq 0$$

$$y_i = 1 : \qquad \alpha(x_i) \geq 1 - \mathcal{G}_i \quad (green)$$

$$y_i = -1 : \qquad \alpha(x_i) \leq +1 + \mathcal{G}_i \quad (red)$$

$$\alpha(x_i) = \theta_0 + \theta_1 x_i$$

$$\left( \hat{\theta}_0, \hat{\theta}_1, \hat{\mathcal{G}} \right) = \boxed{\begin{array}{c} \underset{\theta_0, \theta_1, \mathcal{G}}{argmin} \left\| \theta_1 \right\|^2 + C \sum_{i=1}^{N} \mathcal{G}_i \\[2mm] s.t. \qquad \mathcal{G}_i \geq 0 \\[2mm] y_i \alpha_i(x_i) \geq 1 - \mathcal{G}_i \end{array}}$$

MMC.

CONVEX.

# SVM as loss minimization

$$(\hat{\theta}_0, \hat{\underline{\theta}}_1) = \underset{\theta_0, \underline{\theta}_1}{\mathsf{argmin}} \left( \sum_{i=1}^{N} L(y_i, \alpha(x_i)) \ + \ \lambda \, \|\underline{\theta}_1\|^2 \right)$$

with

$$L(y_i, \hat{y}_i) = \mathsf{max}(0, 1 - y_i \, \hat{y}_i)$$

$$\alpha(x_i) = \theta_0 + x_i \underline{\theta}_1$$



zero-one loss

hinge loss

incorrect prediction

correct prediction

Decision boundary

$$\alpha(x) = 0$$

# Example: SVM on spiral dataset



```
svm = SVC(kernel='linear')
```

```
svm = SVC(kernel='poly',degree=8)
```

```
svm = SVC(kernel='rbf',gamma=1)
```

# Tuning the RBF kernel.



$\gamma = 18.7$

Validation accuracy

$\gamma$