

Kernel Methods

So far : we know how to **compute**
linear predictors.

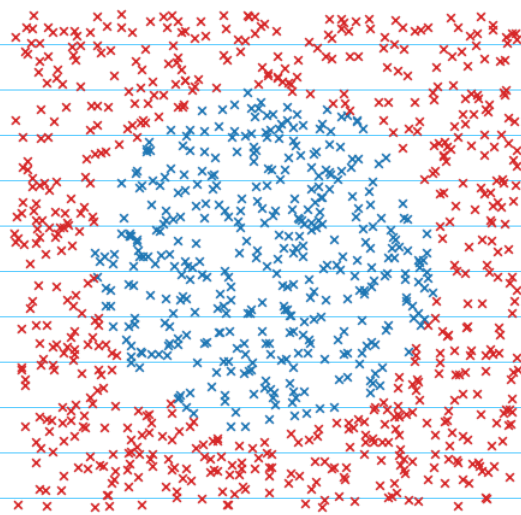
→ Linear Regression

→ Ridge Regression

→ LASSO

→ Logistic Regression

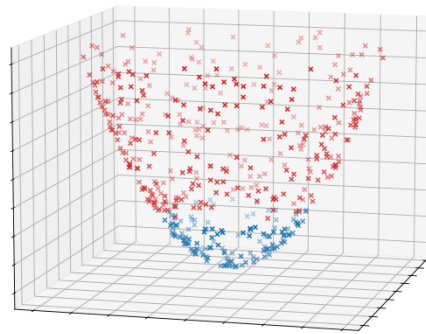
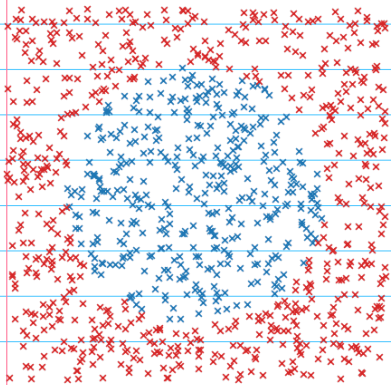
Closed form
or
gradient descent.



① Feature maps:

"If a linear predictor will not work, then lift the dataset to higher dimension and apply a linear method on the lifted dataset."

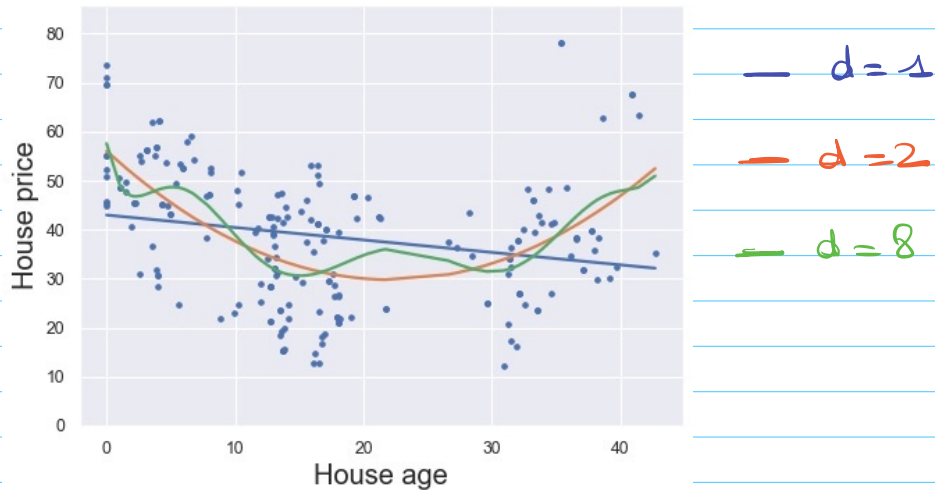
$(x_1, y_1) \dots (x_n, y_n) \quad x_i \in \mathbb{R}^2$



$$x \in \mathbb{R}^2 \xrightarrow{\Phi} \Phi(x) = (x, \|x\|^2) \in \mathbb{R}^3$$

Linear classification on $(\Phi(x_1), y_1) \dots (\Phi(x_n), y_n)$ has a great performance.

• Revisiting Polynomial Regression:



$$x_1, \dots, x_n \in \mathbb{R}$$

$$y_1, \dots, y_n \in \mathbb{R}$$

→ Linear predictor? 😞

Use a **FEATURE MAP** $\Phi: \mathbb{R} \rightarrow \mathbb{R}^{d+1}$

$$\Phi(x) = (1, x, x^2, \dots, x^d)$$

Let $a = (a_0, a_1, \dots, a_d) \in \mathbb{R}^{d+1}$

$$\begin{aligned} \text{Then } \langle a, \Phi(x) \rangle &= a_0 + a_1 x + \dots + a_d x^d \\ &= P_a(x) \end{aligned}$$

→ Polynomial regression aims at minimizing:

$$\begin{aligned} a \in \mathbb{R}^{d+1} &\mapsto \frac{1}{n} \sum_{i=1}^n |y_i - P_a(x_i)|^2 \\ &= \frac{1}{n} \sum_{i=1}^n |y_i - \langle a, \Phi(x_i) \rangle|^2 \\ &= \frac{1}{n} \|Y - Xa\|^2 \end{aligned}$$

$$Y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \quad X = \begin{bmatrix} \Phi(x_1) \\ \vdots \\ \Phi(x_n) \end{bmatrix} = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^d \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^d \end{bmatrix}$$

$\begin{matrix} \uparrow & & \uparrow & & \uparrow \\ & d & & & \\ \leftarrow & & \leftarrow & & \leftarrow \end{matrix}$

→ Polynomial regression
↔ Linear regression on
 $\begin{matrix} \Phi(x_1) & \dots & \Phi(x_n) \\ y_1 & \dots & y_n \end{matrix}$.

- Once again ...
- ① Lift on higher dimension using a feature map
 - ② Apply a linear technique on the transformed dataset.

More Generally: \mathcal{X} general set

$$Y = \mathbb{R}$$

$\Phi: \mathcal{X} \rightarrow \mathbb{R}^D$ feature map

$(x_1, y_1), \dots, (x_n, y_n)$ training sample.

Linear regression on $(\Phi(x_1), y_1), \dots, (\Phi(x_n), y_n)$

$$Q_n: a \in \mathbb{R}^D \mapsto \frac{1}{n} \sum_{i=1}^n |y_i - \langle a, \Phi(x_i) \rangle|^2$$

$$= \frac{1}{n} \|y - \Phi a\|^2$$

$$\Phi = \begin{pmatrix} \Phi(x_1) \\ \vdots \\ \Phi(x_n) \end{pmatrix} \quad n \times D.$$

$$\hat{a} = \underbrace{(\Phi^T \Phi)^{-1}}_{D \times D} \Phi^T y$$

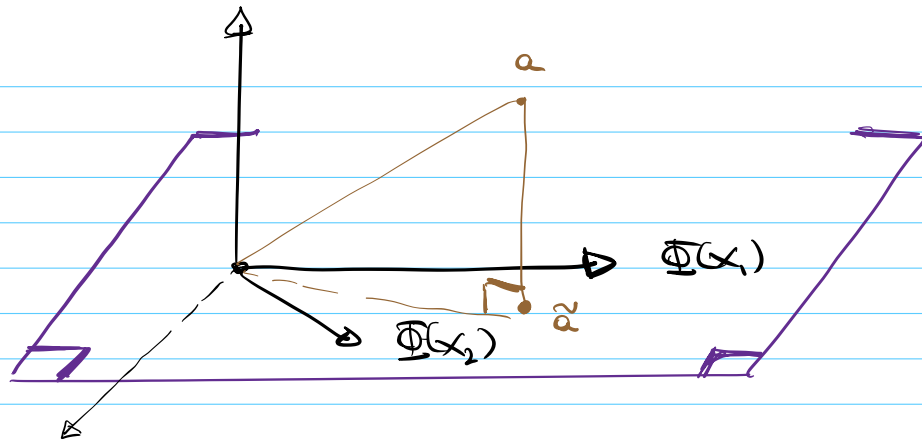
* if $n \geq D$
compute pseudo inverse

Computation cost? $O(D^3)$

What if $D \gg 1$? 😞

⇒ The Kernel Trick

Let \hat{a} be the orthogonal projection of a on $\text{Span}(\Phi(x_1), \dots, \Phi(x_n)) = E$



$$\leadsto \langle a, \Phi(x_i) \rangle = \langle a_2, \Phi(x_i) \rangle$$

$$\Rightarrow \forall a = \forall a_2.$$

We may consider only vectors $a_2 \in E$.

$$\Rightarrow a_2 = \sum_{j=1}^n b_j \Phi(x_j) \text{ for some } b_1, \dots, b_n \in \mathbb{R}.$$

$$\frac{1}{n} \sum_{i=1}^n |y_i - \langle \Phi(x_i), a_2 \rangle|^2$$

$$= \frac{1}{n} \sum_{i=1}^n |y_i - \langle \Phi(x_i), \sum_{j=1}^n b_j \Phi(x_j) \rangle|^2$$

$$= \frac{1}{n} \sum_{i=1}^n |y_i - \sum_{j=1}^n b_j \langle \Phi(x_i), \Phi(x_j) \rangle|^2$$

$$= \frac{1}{n} \|y - Gb\|^2 \quad b = \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix} \in \mathbb{R}^n$$

where

$$G = \begin{bmatrix} \langle \Phi(x_1), \Phi(x_1) \rangle & \dots & \langle \Phi(x_1), \Phi(x_n) \rangle \\ \vdots & & \vdots \\ \langle \Phi(x_n), \Phi(x_1) \rangle & \dots & \langle \Phi(x_n), \Phi(x_n) \rangle \end{bmatrix}$$

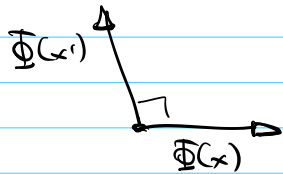
GRAM MATRIX
 $n \times n$

no if G invertible, $\hat{\beta} = G^{-1}y$

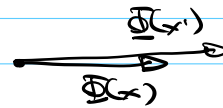
cost $O(n^3)$. $\ll O(D^3)$ if $n \ll D$ || Lift in very high dim

② Reproducing Kernel Hilbert Spaces

The Gram matrix G measures the proximity between the observations.



x and x'
not similar



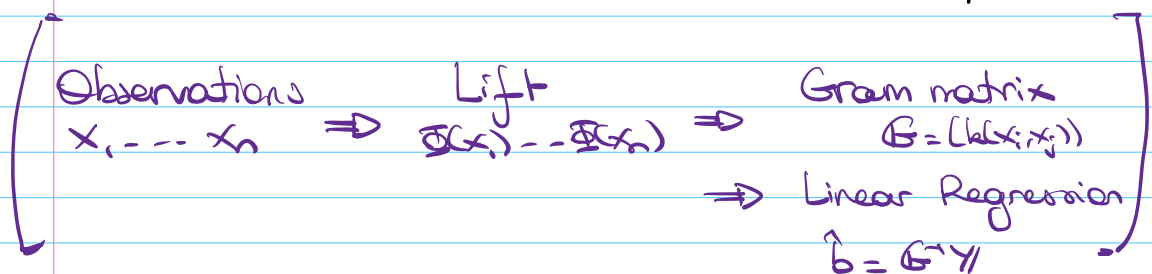
x and x'
similar

→ What if we replace $\langle \Phi(x_i), \Phi(x_j) \rangle$
by a general "measure of similarity"
 $k(x_i, x_j)$.

Ex: RBF / Gaussian kernel

$$k_G(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right)$$

If $k(x, x') = \langle \Phi(x), \Phi(x') \rangle$ for some feature map Φ , then the previous discussion applies.



→ How can we know that k can be written in this way?

• Assume $k(x, x') = \langle \Phi(x), \Phi(x') \rangle$.

Then, $\forall d_1, \dots, d_n \in \mathbb{R}, x_1, \dots, x_n \in \mathcal{X}$,

$$0 \leq \left\| \sum_{i=1}^n d_i \Phi(x_i) \right\|^2 = \sum_{i=1}^n \sum_{j=1}^n d_i d_j \langle \Phi(x_i), \Phi(x_j) \rangle$$

$$\left[0 \leq \sum_{i=1}^n \sum_{j=1}^n d_i d_j k(x_i, x_j) \right] (*)$$

→ If there exists feature map, then

k satisfies $(*) \forall d_1, \dots, d_n \in \mathbb{R}, x_1, \dots, x_n \in \mathcal{X}$

→ This condition is also sufficient!

Hilbert spaces: \mathcal{H} vector space
with a **DOT PRODUCT**

① Symmetry: $\forall x, y \in \mathcal{H}, \langle x, y \rangle_{\mathcal{H}} = \langle y, x \rangle_{\mathcal{H}}$

② Linearity: $\forall x, y, z \in \mathcal{H}, \lambda, \mu \in \mathbb{R}$

$$\langle x, \lambda y + \mu z \rangle_{\mathcal{H}} = \lambda \langle x, y \rangle_{\mathcal{H}} + \mu \langle x, z \rangle_{\mathcal{H}}$$

③ Positive definiteness: $\forall x \in \mathcal{H}, \langle x, x \rangle_{\mathcal{H}} \geq 0$
with $= 0$ iff $x = 0$. $\|x\|_{\mathcal{H}} = \sqrt{\langle x, x \rangle_{\mathcal{H}}}$ is a norm.

④ Completeness: For every continuous linear map $L: \mathcal{H} \rightarrow \mathbb{R}, \exists h \in \mathcal{H}$ such that
 $\forall x \in \mathcal{H}, L(x) = \langle h, x \rangle$.

→ ensures that infinite sums are well-defined - projections

^u Hilbert space \approx Like \mathbb{R}^d but possibly of infinite dimension... ^u

Example: $L_2(\mathbb{R}) = \{f: \mathbb{R} \rightarrow \mathbb{R}, \int |f(x)|^2 dx < \infty\}$

$$\langle f, g \rangle_{L_2(\mathbb{R})} = \int f(x)g(x) dx$$

Def: Let $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$. We say that k is a (positive definite) kernel \Leftrightarrow

① k is symmetric: $\forall x, x' \in \mathcal{X}, k(x, x') = k(x', x)$

② $\forall d_1, \dots, d_n \in \mathbb{R}, x_1, \dots, x_n \in \mathcal{X}$

$$\sum_{1 \leq i, j \leq n} d_i d_j k(x_i, x_j) \geq 0.$$

THEOREM: If k is a kernel, then

there exists $\left\{ \begin{array}{l} \text{Hilbert space } \mathcal{H} \\ \Phi: \mathcal{X} \rightarrow \mathcal{H} \end{array} \right.$ RKHS
Feature map

such that $\forall x, x' \in \mathcal{X}$,

$$k(x, x') = \langle \Phi(x), \Phi(x') \rangle_{\mathcal{H}}.$$

proof sketch:

$$\mathcal{H}_0 = \left\{ \sum_{i=1}^n d_i k(x_i, \cdot); n \in \mathbb{N}, x_i \in \mathcal{X}, d_i \in \mathbb{R} \right\}$$

\hookrightarrow vector space.

$$\begin{aligned} \Phi: \mathcal{X} &\rightarrow \mathcal{H}_0 \\ x &\mapsto k(x, \cdot) \end{aligned}$$

Define $\left\langle \sum_{i=1}^n d_i k(x_i, \cdot), \sum_{j=1}^m d_j k(x_j, \cdot) \right\rangle_{\mathcal{H}_0}$ } This is a dot product

$$:= \sum_{i=1}^n \sum_{j=1}^m d_i d_j k(x_i, x_j)$$

By construction: $\langle \Phi(x), \Phi(x') \rangle_{\mathcal{H}_0}$
 $= \langle k(x, \cdot), k(x', \cdot) \rangle_{\mathcal{H}_0} = k(x, x')$.

Pb: \mathcal{H}_0 is not complete.

It can be completed using a process called completion. \square

How to construct kernels?

k_1, k_2 kernels on \mathcal{X} .

- ① $k_1 + k_2$ is a kernel.
- ② $k_1 \cdot k_2$ is a kernel.
- ③ $\mathcal{X} = \mathbb{R}^d$ $k(x, x') = k(x - x')$

k is a kernel if the Fourier transform

$$\mathcal{F}[k](\xi) = \int e^{-2\pi i \langle x, \xi \rangle} k(x) dx \geq 0.$$

proof:

①

② See Lecture notes

③ $k(x) = \int e^{2\pi i \langle x, \xi \rangle} \mathcal{F}[k](\xi) d\xi$

$$\begin{aligned} \sum_{ij} d_i d_j k(x_i - x_j) &= \int \sum_{ij} d_i d_j \underbrace{e^{2\pi i \langle x_i - x_j, \xi \rangle}}_{\substack{e^{2\pi i \langle x_i, \xi \rangle} \\ e^{-2\pi i \langle x_j, \xi \rangle}}} \mathcal{F}[k](\xi) d\xi \\ &= \langle \sum z_i, \sum \bar{z}_j \rangle = \left\| \sum z_i \right\|^2 \\ &\quad \text{where } z_i = d_i e^{2\pi i \langle x_i, \xi \rangle} \end{aligned}$$

Examples:

① $\Phi: \mathcal{X} \rightarrow \mathcal{H}$ any map. $k(x, x') = \langle \Phi(x), \Phi(x') \rangle_{\mathcal{H}}$ is a kernel.

② $k(x, x') = \langle x, x' \rangle^\alpha \quad \alpha \in \mathbb{N}$

③ Radial Basis Function (RBF) kernel

$$k_{\text{RBF}}(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right)$$

$$\leadsto F[k](f) = \sqrt{2\pi\sigma^2} e^{-\frac{2\pi\sigma^2 \|f\|^2}{\lambda}} \geq 0$$

⑧ Kernel Ridge Regression

Classic Ridge Regression:

$$\begin{array}{l} x_1, \dots, x_n \in \mathbb{R}^d \\ y_1, \dots, y_n \in \mathbb{R} \end{array} \quad \text{Regularization term}$$

$$\text{Minimize } \beta \in \mathbb{R}^d \mapsto \frac{1}{n} \|X\beta - y\|^2 + \lambda \|\beta\|^2$$

$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

$$\leadsto \left[\hat{\beta} = (X^T X + \lambda n I_d)^{-1} X^T y \right]$$

Let's kernelize! \rightarrow any set

$$\begin{array}{l} x_1, \dots, x_n \in \mathcal{X} \\ y_1, \dots, y_n \in \mathbb{R} \end{array}$$

k kernel on \mathcal{X} .

Φ feature map $\mathcal{X} \rightarrow \mathcal{H}$

Kernel Ridge Regression:

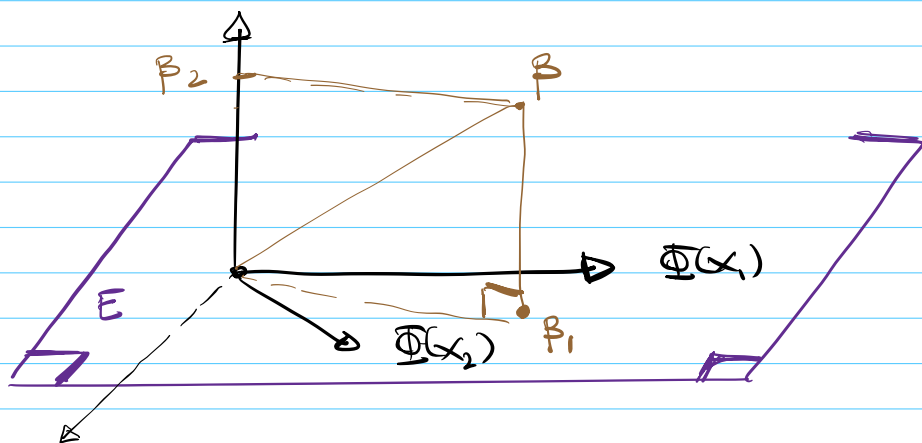
$$(*) \left[\text{Minimize } \beta \in \mathcal{H} \mapsto \frac{1}{n} \sum_{i=1}^n |\langle \Phi(x_i), \beta \rangle - y_i|^2 + \lambda \|\beta\|_{\mathcal{H}}^2 \right]$$

How can we find the minimum?

Representer Theorem:

The minimum of (*) is attained at $\beta \in \mathcal{H}$ of the form $\sum_{i=1}^n a_i \Phi(x_i)$.

proof: Let $E = \text{Span}(\Phi(x_1), \dots, \Phi(x_n))$



$$\beta = \underbrace{\beta_1}_{\in E} + \underbrace{\beta_2}_{\in E^\perp} \quad (= \text{orthogonal of } E)$$

$$\leadsto \langle \beta, \Phi(x_i) \rangle = \langle \beta_1, \Phi(x_i) \rangle$$

$$\leadsto \|\beta\|_{\mathcal{H}}^2 = \|\beta_1\|_{\mathcal{H}}^2 + \|\beta_2\|_{\mathcal{H}}^2$$

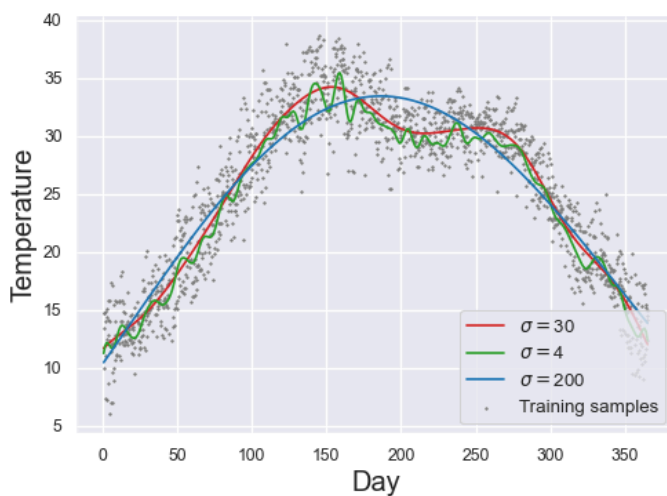
$$\rightarrow \frac{1}{n} \sum_{i=1}^n |\langle \Phi(x_i), \beta \rangle - y_i|^2 + d \|\beta\|_{\mathcal{H}}^2$$

$$= \frac{1}{n} \sum_{i=1}^n |\langle \Phi(x_i), \beta_1 \rangle - y_i|^2 + d \|\beta_1\|_{\mathcal{H}}^2 + \underbrace{d \|\beta_2\|_{\mathcal{H}}^2}_{\geq 0} \quad \square$$

no We minimize over $a = (a_1, \dots, a_n) \in \mathbb{R}^n$
 $\beta_a = \sum_{j=1}^n a_j \Phi(x_j)$

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n |\langle \Phi(x_i), \beta_a \rangle - y_i|^2 + d \|\beta_a\|_2^2 \\ &= \frac{1}{n} \sum_{i=1}^n \left| \sum_{j=1}^n a_j \underbrace{\langle \Phi(x_i), \Phi(x_j) \rangle}_{k(x_i, x_j)} - y_i \right|^2 + d \sum_{i,j} a_i a_j \langle \Phi(x_i), \Phi(x_j) \rangle \\ &= \frac{1}{n} \sum_{i=1}^n \left| \underbrace{G_i^T a}_{= \|Ga - y\|^2} - y_i \right|^2 + d a^T G a \end{aligned} \quad G = (k(x_i, x_j))$$

→ Minimizer $\hat{a} = (G + dn I_n)^{-1} y$.



$x_i = \text{day}$

$$k_g(x, x') = e^{-\frac{(x-x')^2}{2\sigma^2}}$$

↓
Compute G

↓
Compute \hat{a}

⚠ Choice of σ is critical!

$$\sim k_{\sigma}(x, x') \leq 1 \text{ if } \|x - x'\| \leq \sigma \sim$$

Here $\sigma = 30$: \sim if two days are at distance ≤ 30 days, they should be treated similarly. \sim

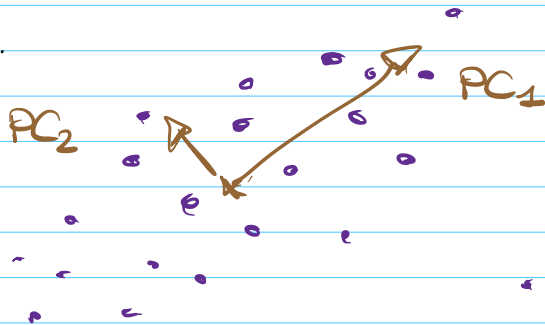
Computations? $\hat{\alpha} = (\underbrace{G + dn I_n}_{n \times n \text{ matrix}})^{-1} y$.

$\leadsto O(n^3)$ = a lot!! (but some tricks, see HW)

Kernel methods are tractable for moderate n :
 $n \lesssim 30000$.

④ Kernel PCA:

Classic PCA:



$$x_1, \dots, x_n \in \mathbb{R}^d \quad X = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$

Principal components = first k eigenvectors
of $XX^T = G$
 $\rightarrow n \times n$ matrix.

\rightarrow Best k -dimensional **LINEAR**
fit of x_1, \dots, x_n
 \leadsto Dimension Reduction

\rightarrow What if no good LINEAR fit?

\rightarrow What if $x_i \notin \mathbb{R}^d$

Kernel PCA k kernel on \mathcal{X}
 $\hookrightarrow \Phi: \mathcal{X} \rightarrow \mathcal{H}$ feature map.

\Rightarrow Apply PCA on $\Phi(x_1), \dots, \Phi(x_n)$.

\leadsto First k eigenvectors of

$$G = (\langle \Phi(x_i), \Phi(x_j) \rangle)_{i,j} \\ = (k(x_i, x_j))_{i,j}$$

Example: Word2Vec: \mathcal{X} = set of words

$$\Phi: \mathcal{X} \rightarrow \mathbb{R}^{25} \quad k(x, x') = \langle \Phi(x), \Phi(x') \rangle$$

x_1, x_2, \dots, x_n = words

↳ either a country or an emotion

germany
kenya

anxious
joy

