# REPRESENTER THEOREM

### Dr. Matthieu R Bloch

Wednesday October 13, 2021

# LOGISTICS

**Assignment 4** due **October 14, 2021**

- Hard deadline on October 16

- Small update posted

**Kayla's office hours tomorrow Thursday October 14, 2021: 11am**

**Assignment 2 grades released:** (2.6/2.7/2.8 not graded)

- Mean: 22.64 - Median: 23.1 - Min: 7.5 - Max: 24.6 (clipped at 24)

**Assignment 3:** 45% graded

**Midterm 1:** 75% graded

**Last time**: solving least squares

- Minimum $\|\cdot\|_2$ solution
- Regularized least squares

Recall:

$$\min_{\theta} \|y - X\theta\|_2^2 + \lambda \|\theta\|_2^2$$

<span style="color:red">hyperparameter</span>

$\lambda > 0$

$X \in \mathbb{R}^{n \times d}$

$$X = \begin{pmatrix} - x_1^T - \\ \vdots \\ - x_n^T - \end{pmatrix} \Big| n$$

$\overset{d}{\longleftrightarrow}$

$$\hat{\theta} = (X^T X + \lambda I)^{-1} X^T y = X^T (X X^T + \lambda I)^{-1} y$$

$\underbrace{\phantom{X^TX}}_{d \times d}$   $\underbrace{\phantom{XX^T}}_{n \times n}$

$$X X^T = \begin{bmatrix} & & \\ & x_i^T x_j & \\ & & \end{bmatrix}_{i,j \in [\![1,n]\!]}$$

# WHAT'S ON THE AGENDA FOR TODAY?

**Last time**: solving least squares
- Minimum $\| \cdot \|_2$ solution
- Regularized least squares

**Today**
- Extension to infinite dimension
- Representer theorem

**Reading:** Romberg, lecture notes 8/9

We can adapt the regularization approach to the situation of a finite dimension Hilbert space $\mathcal{F}$

$$\min_{f \in \mathcal{F}} \sum_{i=1}^{n} (y_i - f(\mathbf{x}_i))^2 + \lambda \|f\|_{\mathcal{F}}^2$$

$\lambda > 0$ hyperparameter

$\in \mathbb{R}$   $\in \mathbb{R}^d$

controls complexity of solution

We can adapt the regularization approach to the situation of a finite dimension Hilbert space $\mathcal{F}$

$$\min_{f \in \mathcal{F}} \sum_{i=1}^{n} (y_i - f(\mathbf{x}_i))^2 + \lambda \|f\|_{\mathcal{F}}^2 \quad (*)$$

*equivalent*

- We are penalizing the norm of the entire function $f$

Using a <u>basis</u> for the space $\{\psi_i\}_{i=1}^{d}$ , and constructing $\boldsymbol{\Psi}$ as earlier, we obtain

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\mathbf{y} - \boldsymbol{\Psi}\boldsymbol{\theta}\|_2^2 + \lambda \boldsymbol{\theta}^\mathsf{T} \mathbf{G} \boldsymbol{\theta} \quad (**)$$

with $\mathbf{G}$ the Gram matrix for the basis.

*depends on basis*          *depends on basis*

Proof of equivalence : we introduce a basis $\{\psi_i\}_{i=1}^{d}$

$$\forall f \in \mathcal{F}, \text{ we can write } f = \underbrace{\sum_{i=1}^{d} \theta_i \psi_i}_{\in \mathbb{R}} \nearrow \in \mathcal{F}$$

$$\min_{f} \underbrace{\sum_{i=1}^{n} |y_i - f(x_i)|^2}_{\textcircled{1}} + \underbrace{\lambda \|f\|_{\mathcal{F}}^2}_{\textcircled{2}}$$

$\textcircled{1}$ : equivalent to $\|y - \underbrace{\Psi \theta}_{\in \mathbb{R}^n}\|_2^2$    $\Psi = \begin{bmatrix} \psi_1(x_1) - - - - \psi_d(x_1) \\ \vdots \\ \psi_j(x_i) \end{bmatrix}$

$\textcircled{2}$ : $\|f\|_{\mathcal{F}}^2 = \|\sum_{i=1}^{d} \theta_i \psi_i\|_{\mathcal{F}}^2 = \left\langle \sum_{i=1}^{d} \theta_i \psi_i, \sum_{j=1}^{d} \theta_j \psi_j \right\rangle = \sum_{i=1}^{d} \sum_{j=1}^{d} \theta_i \theta_j \langle \psi_i, \psi_j \rangle \overset{\Delta}{=} \theta^\top G \theta$

where $\theta = \begin{pmatrix} \theta_1 \\ \vdots \\ \theta_d \end{pmatrix}$    $G = \left[ \langle \psi_i, \psi_j \rangle \right]_{i,j \in [\![1;n]\!]}$

We can adapt the regularization approach to the situation of a finite dimension Hilbert space $\mathcal{F}$

$$\min_{f \in \mathcal{F}} \sum_{i=1}^{n} (y_i - f(\mathbf{x}_i))^2 + \lambda \|f\|_{\mathcal{F}}^2$$

- We are penalizing the norm of the entire function $f$

Using a basis for the space $\{\psi_i\}_{i=1}^{d}$ , and constructing $\mathbf{\Psi}$ as earlier, we obtain

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{\Psi}\boldsymbol{\theta}\|_2^2 + \lambda \boldsymbol{\theta}^\mathsf{T} \mathbf{G} \boldsymbol{\theta}$$

with $\mathbf{G}$ the Gram matrix for the basis.

If $\boxed{\mathbf{\Psi}^\mathsf{T}\mathbf{\Psi} + \lambda\mathbf{G}}$ is invertible, we find the solution as

$$\boldsymbol{\theta}^* = (\mathbf{\Psi}^\mathsf{T}\mathbf{\Psi} + \lambda\mathbf{G})^{-1}\mathbf{\Psi}^\mathsf{T}\mathbf{y} \quad \checkmark$$

⚠ $G \neq I$ in general

and we can reconstruct the function as $f(\mathbf{x}) = \sum_{i=1}^{d} \theta_i^* \psi_i(\mathbf{x}).$ ✗

$$M = X^T X + \lambda G$$

$$M^T = X^T X + \lambda G^T = X^T X + \lambda G \quad \text{is symmetric} \textcolor{blue}{\text{ semidefinite positive}}$$

For any $\theta \in \mathbb{R}^d$ $\quad \theta^T M \theta = \theta^T X^T X \theta + \lambda \theta^T G \theta = \|X\theta\|_2^2 + \lambda \left\| \sum_{i=1}^{d} \theta_i \psi_i \right\|_{\mathcal{F}}^2 \geq 0$
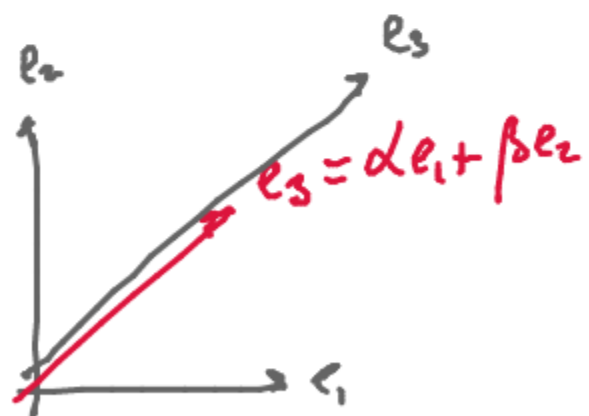
Because $\{\psi_i\}$ is a basis, we know that we can invert $M$

$\left[ \begin{array}{l} \text{Assume } \theta \text{ is such that} \quad \theta^T M \theta = 0 \quad \text{then } \|X\theta\|_2^2 + \lambda \left\| \sum_{r=1}^{d} \theta_i \psi_i \right\|_{\mathcal{F}}^2 = 0 \quad \text{and } \begin{cases} \|X\theta\|_2^2 = 0 \\ \left\| \sum_{i=1}^{d} \theta_i \psi_i \right\|_{\mathcal{F}}^2 = 0 \end{cases} \\ \\ \text{Hence } \theta \in \ker X \text{ and } \sum_{i=1}^{d} \theta_i \psi_i = 0 \quad \text{so that by linear independence } \forall_i \theta_i = 0 \\ \\ \text{Hence } \theta = 0 \end{array} \right.$

Note: $\mathbb{R}^2$



$e_3 = \alpha e_1 + \beta e_2$

We can adapt the regularization approach to the situation of a finite dimension Hilbert space $\mathcal{F}$

$$\min_{f \in \mathcal{F}} \sum_{i=1}^{n} (y_i - f(\mathbf{x}_i))^2 + \lambda \|f\|_{\mathcal{F}}^2$$

- We are penalizing the norm of the entire function $f$

Using a basis for the space $\{\psi_i\}_{i=1}^{d}$ , and constructing $\boldsymbol{\Psi}$ as earlier, we obtain

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\mathbf{y} - \boldsymbol{\Psi}\boldsymbol{\theta}\|_2^2 + \lambda \boldsymbol{\theta}^\mathsf{T} \mathbf{G} \boldsymbol{\theta}$$

with $\mathbf{G}$ the Gram matrix for the basis.

If $\boldsymbol{\Psi}^\mathsf{T}\boldsymbol{\Psi} + \lambda\mathbf{G}$ is invertible, we find the solution as

$$\boldsymbol{\theta}^* = (\boldsymbol{\Psi}^\mathsf{T}\boldsymbol{\Psi} + \lambda\mathbf{G})^{-1}\boldsymbol{\Psi}^\mathsf{T}\mathbf{y} \; = \; \psi^\mathsf{T}\underbrace{\left(\psi\psi^\mathsf{T} + \lambda G\right)^\mathsf{T}}_{\alpha} y \; = \; \text{linear comb o}$$
$$\text{rows of } \psi$$

and we can reconstruct the function as $f(\mathbf{x}) = \sum_{i=1}^{d} \theta_i^* \psi_i(\mathbf{x})$.

If $\mathbf{G}$ is well conditioned, the resulting function is not too sensitive to the choice of the basis

$$\psi_i = \left(\psi_i(x_i) \; \psi_2\right.$$

In $\mathbb{R}^d$, the problem $\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 + \lambda \|\boldsymbol{\theta}\|_2^2$ has a solution

$$\boldsymbol{\theta}^* = \mathbf{X}^\mathsf{T}\boldsymbol{\alpha} \text{ with } \boldsymbol{\alpha} = (\mathbf{X}\mathbf{X}^\mathsf{T} + \lambda\mathbf{I})^{-1}\mathbf{y}$$
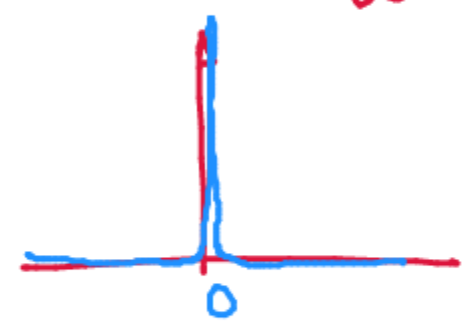
$\mathbf{X}\mathbf{X}^\mathsf{T} \in \mathbb{R}^{n \times n}$ is dimension independent! We will be able to extend this to infinite dimensional Hilbert spaces!

Let $\mathcal{F}$ be a Hilbert space and let $f \in \mathcal{F}$ be the function we are trying to estimate

- We will estimate $f \in \mathcal{F}$ using noisy observations $\boxed{\langle f, x_i \rangle}$ with $\{x_i\}_{i=1}^n$ elements of $\mathcal{F}$

$\hookrightarrow$ for some Hilbert spaces $\langle f, x_i \rangle = f(t_i)$ — fc° of $x_i$

$\triangle$ $f(t) = \int_{-\infty}^{+\infty} f(v) \delta(t-v) dv \leftarrow$

$\int_{-\infty}^{+\infty} \delta(v) dv = 1$

In $\mathbb{R}^d$, the problem $\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 + \lambda\|\boldsymbol{\theta}\|_2^2$ has a solution

$$\boldsymbol{\theta}^* = \mathbf{X}^\mathsf{T}\boldsymbol{\alpha} \text{ with } \boldsymbol{\alpha} = (\mathbf{X}\mathbf{X}^\mathsf{T} + \lambda\mathbf{I})^{-1}\mathbf{y}$$

$\mathbf{X}\mathbf{X}^\mathsf{T} \in \mathbb{R}^{n \times n}$ is dimension independent! We will be able to extend this to infinite dimensional Hilbert spaces!

Let $\mathcal{F}$ be a Hilbert space and let $f \in \mathcal{F}$ be the function we are trying to estimate

- We will estimate $f \in \mathcal{F}$ using noisy observations $\langle f, x_i \rangle$ with $\{x_i\}_{i=1}^n$ elements of $\mathcal{F}$

- This is the equivalent of saying $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{n}$ in finite dimension

**Proposition (Representer theorem)**

$$\min_{f \in \mathcal{F}} \sum_{i=1}^n \left| y_i - \langle f, x_i \rangle_{\mathcal{H}} \right|^2 + \lambda\|f\|_{\mathcal{H}}^2$$

has solution

$$f = \sum_{i=1}^n \alpha_i x_i \text{ with } \boldsymbol{\alpha} = (\mathbf{K} + \lambda\mathbf{I})^{-1}\mathbf{y} \qquad \mathbf{K} = [\langle x_i, x_j \rangle]_{1 \le i,j \le n}$$

$\in \mathcal{F}$

$\in \mathbb{R}^{n \times n}$

$\in \mathbb{R}$

**Proof.**

$$\min_{f} \sum_{i=1}^{n} |y_i - \langle f, x_i \rangle|^2 + \lambda \|f\|^2 \qquad (*)$$

$\forall f$ we can write $f = f_1 + f_2$ w/ $f_1 \in \text{Span}(\{x_i\}_{i=1}^{n})$ and $f_2 \in \text{Span}(\{x_i\}_{i=1}^{n})^{\perp}$

$\qquad\qquad\qquad\qquad\qquad\;\; \hookleftarrow$ closest point to $f$ in $\text{Span}(\{x_i\}_{i=1}^{n})$

Note that $\forall i \;\; \langle f_2, x_i \rangle = 0$

Then $\displaystyle\sum_{i=1}^{n} |y_i - \langle f, x_i \rangle|^2 = \sum_{i=1}^{n} |y_i - \langle f_1, x_i \rangle|^2$

$\|f\|^2 = \|f_1 + f_2\|^2 = \|f_1\|^2 + \underbrace{\|f_2\|^2}_{\geq 0}$ (Pythagorean theorem)

Hence any solution of $(*)$ must live in $\text{Span}(\{x_i\}_{i=1}^{n})$

# LEAST-SQUARES IN INFINITE DIMENSION HILBERT SPACES

In $\mathbb{R}^d$, the problem $\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 + \lambda\|\boldsymbol{\theta}\|_2^2$ has a solution

$$\boldsymbol{\theta}^* = \mathbf{X}^\mathsf{T}\boldsymbol{\alpha} \text{ with } \boldsymbol{\alpha} = (\mathbf{X}\mathbf{X}^\mathsf{T} + \lambda\mathbf{I})^{-1}\mathbf{y}$$

$\mathbf{X}\mathbf{X}^\mathsf{T} \in \mathbb{R}^{n \times n}$ is dimension independent! We will be able to extend this to infinite dimensional Hilbert spaces!

Let $\mathcal{F}$ be a Hilbert space and let $f \in \mathcal{F}$ be the function we are trying to estimate

- We will estimate $f \in \mathcal{F}$ using noisy observations $\langle f, x_i \rangle$ with $\{x_i\}_{i=1}^n$ elements of $\mathcal{F}$

- This is the equivalent of saying $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{n}$ in finite dimension

**Proposition (Representer theorem)**

$$\min_{f \in \mathcal{F}} \sum_{i=1}^n \left| y_i - \langle f, x_i \rangle_{\mathcal{H}} \right|^2 + \lambda\|f\|_{\mathcal{H}}$$

has solution

$$f = \sum_{i=1}^n \alpha_i x_i \text{ with } \boldsymbol{\alpha} = (\mathbf{K} + \lambda\mathbf{I})^{-1}\mathbf{y} \qquad \mathbf{K} = [\langle x_i, x_j \rangle]_{1 \leq i,j \leq n}$$

We will see that the situation of the representer theorem happens in Reproducing Kernel Hilber Space (RKHS)