



Machine Learning and Micromechanics as Allies to Establish Composition-Property correlations in Cement Pastes

Tulio Honorio, Sofiane Ait Hamadouche, Amelie Fau

► To cite this version:

Tulio Honorio, Sofiane Ait Hamadouche, Amelie Fau. Machine Learning and Micromechanics as Allies to Establish Composition-Property correlations in Cement Pastes. 2022. hal-03723418v1

HAL Id: hal-03723418

<https://hal.science/hal-03723418v1>

Preprint submitted on 14 Jul 2022 (v1), last revised 26 Jan 2023 (v3)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

Machine Learning and Micromechanics as Allies to Establish Composition-Property correlations in Cement Pastes

 **Tulio Honorio**¹,  **Sofiane Ait Hamadouche**¹, and  **Amelie Fau**¹

¹ Université Paris-Saclay, CentraleSupélec, ENS Paris-Saclay, CNRS, LMPS - Laboratoire de Mécanique Paris-Saclay, 91190, Gif-sur-Yvette, France

Composition-property correlations are fundamental to understand cement-based materials behavior and optimize their formulation. Modelling based on fundamental material component constitutes a reliable tool to establish these correlations with the advantage of better exploring formulation space when compared with the often adopted experimental trial-and-error approaches. In this context, Machine Learning (ML) and Micromechanics-Based (MB) methods have been concurrently used for property prediction from material composition. Here, we show that these techniques can be allies for establishing composition-property correlations. We focus on predictions of Ordinary Portland Cement pastes elastic properties, but the outlined strategy can be extended to other cement systems. Various microstructures representations are considered in MB estimates, including multiscale representations and representations with ellipsoidal inclusions. In contrast, ML predictions do not need any a priori assumption on material microstructure. Predictions using ML and MB yield similar accuracy when compared against test datasets (but ML performed much better regarding the error estimated in training datasets). Working as allies, ML can be deployed to evaluate the (lack of) knowledge over the multi-dimensional parametric domains, and micromechanics provides a theoretical background for property data curation and is a tool to make up for missing data in databases.

Keywords machine learning; micromechanics; ordinary portland cement paste; elastic properties; data science; early ages.

1 Introduction

Establishing Processing - Composition - (Micro) Structure - Property - Performance correlations is the central paradigm for understanding material behavior in a bottom-up perspective as well as to conceive and optimize materials for tailored applications (Olson 1997). Such correlations are important for cement-based materials since, on the *composition* side, the key ingredients vary largely according to the local availability of resources, *processing* spans lower and higher technology contexts (Wangler et al. 2019), and the design of cement components and concrete structures relies on *property and performance* requirements. Material property prediction having as input the composition is therefore critical to optimize the use of cement-based materials.

Micromechanics-based (MB) modeling have been successfully used to unveil Composition - Property correlations for various properties in cement-based materials, including mechanical (Wyrzykowski et al. 2017; Pichler and Hellmich 2011; Sanahuja et al. 2007; Königsberger, Honório, et al. 2021), transport and thermal (Bary et al. 2006; Patel et al. 2016; Honorio, Bary, and Benboudjema 2018), and electromagnetic (Guihard et al. 2019; Honorio, Carasek, et al. 2020; Honorio, Bore, et al. 2020) properties, as well as coupling properties in the thermo-poro-mechanical framework (Ulm et al. 2004; Ghabezloo et al. 2009; Honorio, Bary, and Benboudjema 2018; Honorio, Brochard, et al. 2018). An advantage of MB modeling is the simplicity of computations, which enables assessing various scenarios of interest regarding the composition, uncertainty on phase properties (Honorio, Carasek, et al. 2020), and morphology of phases in a heterogeneous material. However, one may legitimately dispute the pertinence of representing the microstructure of cement paste under the usual assumptions adopted in analytical homogenization approaches. These assumptions include (i) a random microstructure (there is evidence that some correlation between phases volume distribution has been quantified using microstructural

hydration model (Hlobil 2020)), (ii) phases being often represented by spherical (or ellipsoidal) inclusions (experimental evidence shows that crystalline phases cement paste are not generally spherical or ellipsoidal), (iii) perfect interfaces among phases (while some defects may exist), and (iv) separability of scale (especially considering that heterogeneity size, for example of cement particles, may span various magnitudes). Numerical homogenization is not immune to the same questionings.

In this context, Machine Learning (ML) arise as a promising tool to directly establish Composition - Property correlations without *a priori* assumption of the microstructure characteristics (Agrawal et al. 2016). The huge amount of experimental data produced on cement-based materials in the last century can be used to build databases that can be interrogated by ML. As highlighted by Bullard et al. (Bullard et al. 2017), a “systematic development of structure-property relationships” based on both the “curation of fundamental material component data” and “validated modeling based on fundamental scientific principles” may “revolutionize” the design of cement-based materials. However, as recognized by the same authors, such an approach was given comparatively little attention in the concrete research community when compared to the “increasingly laborious trial-and-error exploration of the design space and mixture qualification process” (Bullard et al. 2017). In cement-based materials research, ML has been deployed since the 90’s to predict the compressive strength (Kasperkiewicz et al. 1995; I. -. Yeh 1998; I.-C. Yeh et al. 2009; Duan et al. 2013; Young et al. 2019) using frequently artificial neural networks (ANN). Other methods include support vector machines (Yan et al. 2010), decision trees (Behnood et al. 2015), evolutionary algorithms (Golafshani et al. 2018). Elastic properties have been also extensively studied using ML (Ben Chaabene et al. 2020), with a strong focus on the impact of using recycled aggregates. As input variables, the composition in terms of cement and water content, as well as supplementary cementitious materials (SCM) and admixture mass or volume, are often adopted (Ben Chaabene et al. 2020). Neither the effects of the mineralogical composition of cement nor the effects of age (and property development, especially at early-age) are generally considered in these studies.

In this work, a multi-technique modeling approach combining ML and MB methods is proposed to link cement system composition and degree of hydration to the elastic properties of the material. We tackle specifically the predictions of OPC pastes elastic properties from the composition of the cement (in terms of clinker composition and gypsum fraction), w/c and age, but the strategy outlined here can be extended to other cement systems and scales. Since OPC system are simpler and better experimentally characterized than other cement systems, they are an ideal candidate for testing the approach presented here and for demonstrating its feasibility. We explore paths in which ML and MB techniques can be allies, notably in the analysis of experimental databases to evaluate existing experiments and lack of experiments and by providing missing data. The results obtained are a contribution towards the development of multiscale modeling of cement-based materials informed by the cement composition variability and enhanced by blending data from different research projects. This framework can be used to improve the comprehension of correlations among the composition, microstructure, and properties of cement-based materials.

2 Machine Learning approach and database construction for predicting elastic properties

Knowledge about cement paste and behavior is fundamentally offered through experimental observations. A direct approach to exploit the large literature is collecting wide range of experimental results published and using machine learning methods to predict properties for new compositions based on the training dataset.

2.1 A database construction for cement pastes linking composition and elastic constants

Based on experimental data from the literature (Helmuth et al. 1966; Haecker et al. 2005; Boumiz et al. 1996; Tamtsia et al. 2004; Wang et al. 2011; Constantinides et al. 2004; Lura et al. 2003; Chamrova 2010; Sun et al. 2007; Maruyama et al. 2014) a dataset with 376 entries is built, which will be used for training and validation. Details on database construction are given in the

Appendix. The inputs in the datasets are cement composition (in terms of clinker minerals and gypsum contents), water-cement ratio, age, and degree of hydration. The outputs are the elastic constants: E Young, K bulk and G shear moduli and ν Poisson ratio. Note that the dimensionality of the manifold can be reduced considering that the elastic constants are linked through simple relations in the case of isotropic materials ($E = 9KG/(3K + G)$; $\nu = (3K - 2G)/(2(3K + G))$; $K = E/(3(1 - 2\nu))$ and $G = E/(2(1 + \nu))$). Also, the age and the degree of hydration can be linked using a bijection (e.g. a sigmoid function).

Table 1: Statistical analysis of the cement paste dataset of 365 observations used for training. * dimensionless.

Data	Variable	Min.	Max.	Mean	St. Dev.	Exceed Kurtosis*	Skewness*
Input	Age [days]	0.12	720	49	124	13.9	3.6
	DOH [-]	0.03	1.0	0.5	0.3	-0.7	0.6
	w/c [-]	0.25	0.80	0.44	0.10	-0.20	0.64
	m_{C_3S} [%]	24.5	100.0	60.2	13.8	2.1	-0.7
	m_{C_2S} [%]	0.0	61.3	16.6	15.2	2.4	1.6
	m_{C_3A} [%]	0.0	12.7	8.1	3.4	-0.7	-0.6
	m_{C_4AF} [%]	0.0	12.7	5.8	4.2	-1.5	-0.2
	m_{Gypsum} [%]	0.0	6.8	2.9	2.9	-1.7	0.3
Output	E [GPa]	0.22	37.2	11.2	7.8	0.3	0.8
	ν [-]	0.07	0.49	0.30	0.07	0.88	0.58
	K [GPa]	0.15	32.2	9.3	5.6	2.0	1.3
	G [GPa]	0.07	14.6	4.4	3.1	0.3	0.8

Table 1 shows the statistical parameters associated with the training dataset. In the various ML applications for mechanical properties of cement-based materials, the dataset size spans from 74 (Ben Chaabene et al. 2020) up to more than 10,000 (Young et al. 2019) observations, most of the cases with data size in the range 100 to 1,000 observations (Ben Chaabene et al. 2020). The size of the dataset provided here has, therefore, an intermediary size. It can already provide sufficient support for learning but could surely be improved with complementary data in future works.

2.2 Machine Learning methods

For *prediction* purposes the following algorithms are employed:

- **Linear Regression (LR):** The output is predicted using a linear combination of the numerical features vector. The conditional probability is computed using a parameter vector estimated from the minimization of a loss function.
- **Decision Tree (DT):** A decision tree (i.e. a flow chart structure in which the internal nodes correspond to a test on a feature, while the branches correspond to an outcome of the test) is built using Classification and Regression Trees (CART) algorithm (Breiman et al. 1984).
- **Gradient Boosted Trees (GBT):** A prediction model is constructed in the form of an ensemble of trees, which is trained sequentially in order to enhance the capability of the previous trees. The implementation adopted is based on LightGBM algorithm (Ke et al. 2017).
- **Nearest Neighbors (NN):** This instance-based learning technique predicts a value by analyzing the nearest neighbors in the feature space.
- **Artificial Neural Network (ANN):** A neural network is constituted of stacked layers, each associated with simple computation. The information is processed layer by layer starting at the input layer until the output layer. The neural network is trained in order to minimize a loss function on the training set. A gradient descent method is used to perform this minimization.
- **Random Forest (RF):** Various decision trees are constructed and the prediction is done by taking the mean value of the tree predictions based on bootstrap aggregating algorithm

(Breiman 1996), each decision tree is trained using only a random subset of the features.

- **Gaussian Process (GP):** Predictions are made using Bayesian inference on the Gaussian process conditioned to the training data e.g. (Williams et al. 1996). The underlying assumption of the method is that the prediction function can be associated with a Gaussian process (defined by its kernel or covariance function). The training phase consists of estimating the parameters of the kernel.

We use Mathematica 12.0.0.0 software (Inc. n.d.), in which these algorithms are built-in. Methods like ANN, LR, and GP produce a smooth predictor, whilst DT, NN and RF produce discrete prediction values.

2.3 Validation process

The performances of the ML methods are estimated using a k -fold cross-validation technique (Bengio et al. 2004). The training dataset is divided in k folds, i.e subsets \mathcal{D}_i in which elements are randomly sampled from the dataset. In each fold construction, care is taken so that a given element is not chosen more than once (in order to ensure that the intersection set of all folds is the empty set: $\mathcal{D}_i \cap \mathcal{D}_j = \emptyset, \forall (i, j) \in [1; k]^2$ if $i \neq j$). The predictor is trained on $k - 1$ folds and, then, it is used to predict the values in the remaining fold. This operation is repeated k times so that all folds have been used for validation. Here, use k -fold method with $k = 5$ and 10 folds, as usually is done in the literature (Nematzadeh et al. 2015; Rodriguez et al. 2010).

3 Estimation of elastic properties from micromechanics

MB approaches have been proven useful to get accurate predictions (Wyrzykowski et al. 2017) of properties of cement systems at various scales, sometimes with error not even exceeding 3% (Königsberger, Honório, et al. 2021). In the literature, various propositions of representation of the cement paste microstructure exist based on different number of scales, system morphology, and on the use of different models to describe the volume fraction of the constituents in the system. To fully explore the relevant microstructure representations mostly adopted, here we consider sixteen representations (see details in Section 3.2). Each representation combines different assumptions regarding the number of scales to be considered, the shape of the constituent phases or the model to describe phase assemblage. These representations are based on previous studies on the upscaling of different physical properties of cement-based materials (Sanahuja et al. 2007; Honorio, Bary, and Benboudjema 2016; Honorio, Bary, and Benboudjema 2018; Königsberger, Honório, et al. 2021).

Powers (Powers and Brownard 1946), Königsberger-Hellmich-Pichler (KHP) (Königsberger, Hellmich, et al. 2016), and Tennis and Jennings (Tennis et al. 2000) models are considered to evaluate the phases evolution with the degree of hydration in OPC pastes. The former is the earliest and ones of the most simple strategies. The latter is one of the most detailed descriptions of phase assemblage in OPC systems before resorting to thermodynamics modelling. KHP model updates Powers model by introducing C-S-H densification. In the following, we detail how we obtain the input for micromechanics estimations (i.e volume fractions of phases as a function of the age or DOH). Then, the formulation of the homogenization schemes are recalled.

3.1 Phase assemblage approximation from hydration models for Ordinary Portland Cement (OPC) pastes

Powers hydration model (Powers and Brownard 1946) considers only three phases, as listed in Table 2. It has been coupled with micromechanics strategies to study early-age property development of cement-based materials in (Sanahuja et al. 2007; Pichler, Hellmich, et al. 2013). This model has the advantage of its simplicity, but it does not account for a variety of phases that can be present in OPC systems.

KHP model extends Powers model by considering C-S-H densification (in agreement with NMR data) and by providing the volume fraction of portlandite.

For comparison, a more elaborate model, the Tennis and Jennings (Tennis et al. 2000) model, is explored. It describes the chemical rearrangement due to the hydration process by stoichiometric relationships based on a more detailed separation of phases i.e. the evolution of clinker minerals

and gypsum fractions as well as the main hydrates separately as a function of the degree of hydration, as listed in Table 2. It even allows to distinguish low-density (LD) and high-density (HD) C-S-H.

More details about formulations of Powers, and Tennis and Jennings models are given in Appendix B.

The elastic properties of the constituent phases are given in Table 2.

Table 2: Elastic constants of phases. * Monosulfoaluminate. ** Dihydrate. ^a Molecular simulations.

Hyd. model	Phase	E [GPa]	ν [GPa]	G [GPa]	K [GPa]	Ref.
Powers/KHP	Clinker	140	0.30	53.8	116.7	Acker (2001)
	Hydrates	22.06	0.24	11.76	18.69	Pichler and Hellmich (2011)
	Pores	0	0.5	0	2.18	Hammond (1997)
TJ	C ₃ S	135 ± 7	0.3	51.9	112.5	Velez et al. (2001)
	C ₂ S	130 ± 20	0.3	50.0	108.3	Velez et al. (2001)
	C ₃ A	145 ± 10	0.3	55.8	120.8	Velez et al. (2001)
	C ₄ AF	125 ± 25	0.3	48.1	104.2	Velez et al. (2001)
	C \bar{S} H ₂ **	45.7	0.33	17.2	44.8	Aller et al. (1996)
	HD C-S-H	29.4 ± 2.4	0.24	30.4	18.8	Constantinides et al. (2004)
	LD C-S-H	21.7 ± 2.2	0.24	11.9	13.9	Constantinides et al. (2004)
	CH	42.0	0.315	16.0	37.8	Monteiro et al. (1995)
	AFt	25.0 ± 2	0.34 ± 0.02	9.3	26.0	Speziale et al. (2008)
	AFm*	24.5	0.34	9.1	25.5	Honorio, Guerra, et al. (2020) ^a
	C ₄ AH ₁₃	25.0	0.34	9.3	26.0	Speziale et al. (2008)
	Hydrogarnet	55.5	0.35	20.6	61.7	Manzano (2009) ^a
	Pores	0	0.5	0	2.18	Hammond (1997)

3.2 Representations of the microstructure

Sixteen representations of the microstructure of the cement paste are considered here, each one combining different assumptions regarding the number of scales to be considered, the shape of the constituent phases or the model to describe phase assemblage.

Figure 1 shows the eight representations of the microstructure of the cement paste tested in the case of spherical inclusions. The other eight representations refer to the adoption of ellipsoidal inclusions to represent some phases. For the cases with ellipsoidal inclusions, similar representation are adopted with the following modifications: (i) C-S-H (when TJ model is used) or hydrates (when Powers model is used) are modelled as a elongated inclusions with aspect ratio of 10; and (ii) AF-phases and CH (when TJ model is used) are modelled as oblate particles with aspect ratio of 0.2. All the other phases including pores are considered as spherical inclusions.

Using the description of phases by the *Powers model*, the Mori-Tanaka, with hydrates functioning as the matrix, and Self-Consistent schemes are concurrently considered, which leads to the MTPow and SCPow macroscopic behaviors, respectively.

Using the Tennis and Jennings model, in addition to the flexibility offers by the two upscaling schemes, microstructure can be constructed with different perspectives. All hydrates, anhydrides, and pores can be treated at the same scale, which gives MT1s and SC1s corresponding with Mori-Tanaka with LD C-S-H as matrix and self-consistent schemes, respectively. Or, C-S-H gel can be handled at a smaller scale comprising LD and HD C-S-H domains and gel porosity. First, the effective properties of C-S-H gel are obtained using SC scheme on a heterogeneous material. The effective properties of C-S-H gel are then used in parallel with the properties of other hydrates and clinker inclusions at the cement paste scale as input for the second stage of homogenization, which can be processed using Mori-Tanaka with C-S-H gel as the hosting matrix or self-consistent schemes to give MT2s and SC2s effective properties.

Using the description of phases by the KHP model, a two-scale representation is considered with C-S-H gel scale and a cement paste scale *per se* at the higher level.

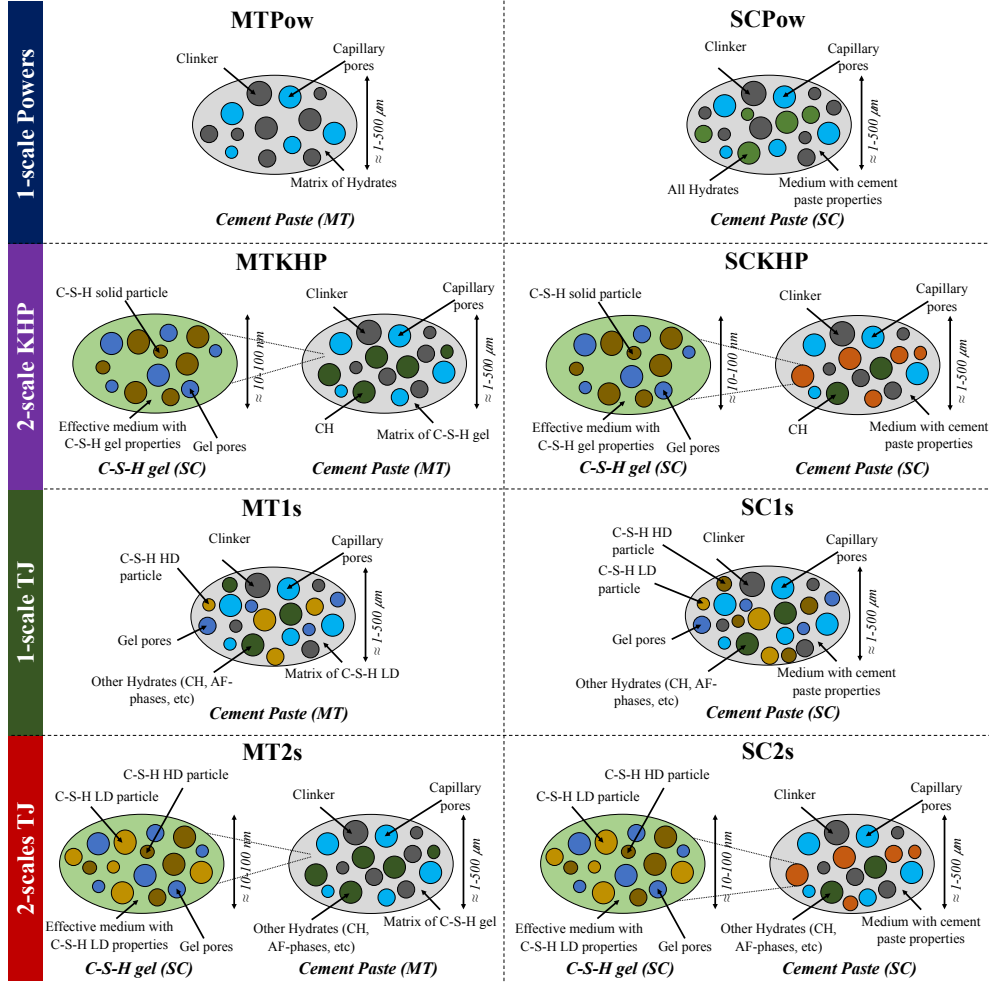


Figure 1: Cement paste microstructure considered for micromechanics upscaling schemes: Input volume fractions are obtained from Powers or Tennis and Jennings (TJ) models. SC or MT are deployed to upscale cement paste elastic properties. For 2-scale representations, C-S-H gel effective properties are upscaled using SC scheme. All the other phases including pores are considered as spherical inclusions. SC scheme is deployed to upscale cement paste elastic properties. For the representations with ellipsoidal inclusions: with Powers model, hydrates are considered as prolate particles $a_r = 10$; With TJ model, C-S-H LD and HD are considered as prolate particles $a_r = 10$, and CH and AFm as oblate particles $a_r = 0.2$

3.3 Analytical homogenization of the elastic properties of micro and macro isotropic heterogeneous materials

We deploy Mori-Tanaka (MT) and Self-Consistent (SC) homogenization scheme for micro and macro-isotropic heterogeneous materials with ellipsoidal inclusions randomly distributed in a representative elementary volume (REV). According to these schemes, the effective stiffness tensor C_{est} , with the superscript *est* designating MT or SC estimate, of a heterogeneous material is given by (e.g. (Zaoui 2002)):

$$C_{est} = \left(\sum_{r=1}^N f_r C_r : [I + P^0 : (C_r - C^0)]^{-1} \right) : \left(\sum_{r=1}^N f_r [I + P^0 : (C_r - C^0)]^{-1} \right)^{-1} \quad (1)$$

where f_r is the volume fraction of the phase r , C_r is the stiffness tensor of phase r ; $P^0 = S_H^0 : C^0$ is the Hill tensor obtained from the Eshelby tensor S_H^0 (which depends only on the properties of the reference medium, see ref. (Mura 1987) for the expressions of Eshelby tensors including the case of ellipsoidal inclusions) and the stiffness tensor of the reference medium C_0 , which is defined according to the scheme chosen:

- $C^0 = C_0$ for the MT scheme, where the subscript 0 stands for matrix properties.
- $C^0 = C^{SC}$ for the SC scheme, i.e. the reference medium is the effective medium itself.

An important input for estimations using non-spherical particles is the aspect ratio of the

particles. We adopt an aspect ratio of $a_r = 10$ (prolate particle) for C-S-H needles and $a_r = 0.2$ (oblate particle) for crystalline hydrates such as CH and AFm.

In the case of spherical isotropic inclusions, Eq. 1 simplifies into the forms described below. For an $(N + 1)$ -phase heterogeneous material with a matrix/inclusion morphology constituted of N isotropic spherical inclusions randomly distributed in a matrix (percolating phase), the Mori-Tanaka estimate of the effective bulk K^{MT} and shear G^{MT} moduli are, respectively, obtained from (e.g. (Torquato 2002)):

$$\frac{K^{MT} - K_0}{K^{MT} + \frac{4}{3}G_0} = \sum_{r=1}^N f_r \frac{K_r - K_0}{K_r + \frac{4}{3}G_0}; \quad \frac{G^{MT} - G_0}{G^{MT} + \frac{4}{3}H_0} = \sum_{r=1}^N f_r \frac{G_r - G_0}{G_r + \frac{4}{3}H_0} \quad (2)$$

with $H_0 = \frac{\frac{3}{2}K_r + \frac{4}{3}G_r}{K_r + 2G_r} G_r$, and the subscript 0 denotes the (isotropic) matrix phase.

For an N -phase heterogeneous materials with N isotropic equiaxed inclusions randomly distributed in representative elementary volume following a polycrystalline-like morphology (i.e. in which no phase clearly functions as a matrix), the Self-Consistent effective bulk K^{SC} and shear G^{SC} moduli are given, respectively, by the implicit relations (e.g. (Torquato 2002)):

$$\sum_{r=1}^N f_r \frac{K_r - K^{SC}}{K_r + \frac{4}{3}G^{SC}} = 0; \quad \sum_{r=1}^N f_r \frac{G_r - G^{SC}}{G_r + H_{SC}} = 0. \quad (3)$$

3.4 Bounds for the elastic properties

From the properties of the constituent phases and their volume fraction, micromechanics offers not only the effective properties but also bounds between which the elastic properties of the heterogeneous material should lie within. It is then possible to cross-check the observed experimental values with the bounds given by the theoretical models based on specific modelling assumptions.

Two theoretical bounds defined in terms of the effective bulk K^{eff} and shear G^{eff} moduli are considered in this paper (Zaoui 2002):

- **Voigt-Reuss** bounds, which are associated to series and parallel models:

$$\left(\sum_{r=1}^N \frac{f_r}{K_r} \right)^{-1} \leq K^{eff} \leq \sum_{r=1}^N f_r K_r; \quad \left(\sum_{r=1}^N \frac{f_r}{G_r} \right)^{-1} \leq G^{eff} \leq \sum_{r=1}^N f_r G_r \quad (4)$$

where the leftmost term is the Reuss estimate and the rightmost term, the Voigt estimate.

- **Hashin-Shtrikman** (HS) bounds are defined for heterogeneous materials with an isotropic distribution of phases for an arbitrary phase geometry based on variational principle in linear elasticity (Hashin et al. 1963):

$$\sum_{r=1}^N \frac{f_r K_r / [K^- + \alpha^-(K_r - K^-)]}{f_r / [K^- + \alpha^-(K_r - K^-)]} \leq K^{eff} \leq \sum_{r=1}^N \frac{f_r K_r / [K^+ + \alpha^+(K_r - K^+)]}{f_r / [K^+ + \alpha^+(K_r - K^+)]} \quad (5)$$

$$\sum_{r=1}^N \frac{f_r G_r / [G^- + \beta^-(G_r - G^-)]}{f_r / [G^- + \beta^-(G_r - G^-)]} \leq G^{eff} \leq \sum_{r=1}^N \frac{f_r G_r / [G^+ + \beta^+(G_r - G^+)]}{f_r / [G^+ + \beta^+(G_r - G^+)]} \quad (6)$$

where $G^- = \inf(G_r)$; $K^- = \inf(K_r)$; $G^+ = \sup(G_r)$; $K^+ = \sup(K_r)$ are the extreme values of the bulk and shear moduli considering all r phases; $\beta^\pm = \frac{6(K^\pm + 2G^\pm)}{5(3K^\pm + 4G^\pm)}$ and $\alpha^\pm = \frac{3K^\pm}{3K^\pm + 4G^\pm}$. HS bounds are narrower than Voigt-Reuss bounds.

The bounds for the Young modulus can be directly computed from the lower and upper bounds using (Zimmerman 1992):

$$\frac{9K_L G_L}{3K_L + G_L} \leq E^{eff} \leq \frac{9K_U G_U}{3K_U + G_U} \quad (7)$$

where the subscript L refers to the lower (HS or Reuss) bound; and the subscript U , to the upper (HS or Voigt) bound.

For the Poisson ratio, Zimmerman (Zimmerman 1992) shows that the correct bounds are given by:

$$\frac{3K_L - 2G_U}{6K_L + 2G_U} \leq \nu^{eff} \leq \frac{3K_U - 2G_L}{6K_U + 2G_L} \quad (8)$$

where the largest possible value of ν refers to the largest value of K combined with the smallest value of G , and vice versa. The argument is valid for both Voigt-Reuss and HS bounds.

4 Results and Discussion

MB and ML methods are investigated for predictions and analysis of various properties of cement paste. Then, bounds for elastic properties given by MB methods are compared with experimental observations. Predictions of elastic properties given by MB and ML methods are compared for training and test datasets. Finally, the lack of knowledge on the parametric input is evaluated and the experimental dataset is enriched with MB observations guided by ML evaluations.

4.1 Micromechanics bounds for dataset curation

Knowing w/c and DOH (or age, from which DOH can be estimated), fractions of phases are evaluated from hydration models, and bounds for E , ν , G and K are derived from Voigt-Reuss and Hashin-Shtrikman theories (as detailed in Section 3.4). Comparing the experimental elastic properties and the bounds, both lower bounds, being null, are satisfied by all experimental observations. However, some experimental observations of K and ν exceed the upper bounds. Proportions of values exceeding the theoretical bounds are summarized in Table 3. Since the phase intrinsic properties are associated with a variability/uncertainty on the order of 10-20% as reported in Table 2, we also provide bounds estimation accounting for an average 15% uncertainty.

Table 3: Fraction of values exceeding the upper bounds, Voigt and Hashin-Shtrikman (HS), for each elastic constant tested. Bounds computed using Powers (Pow.) or Tennis and Jennings (TJ) hydration models and considering a 15% uncertainty on K and G reported as phase properties in Table 2. The values in-between the parenthesis refer to bounds computed using the average values reported in Table 2 (i.e. without the 15% uncertainty).

Upper bounds	Hydration model	E	ν	K	G
Voigt	Pow.	0 (0)	0 (0.002)	0 (0.8)	0 (0)
Voigt	TJ	0 (0)	0 (0)	0 (0.5)	0 (0)
HS	Pow.	0 (0)	0 (0)	1.4 (4.1)	0 (0)
HS	TJ	0 (0)	0 (0.008)	1.4 (3.5)	0 (0)

All values of both shear and Young moduli are below the upper bounds. A few values of the bulk modulus, less than 5% for the worst case of the experimental observations, exceed the upper bounds, when the uncertainty on the phase elastic moduli are not accounted for. As expected, more points exceed HS than Voigt bound since the Hashin-Shtrikman bounds are tighter. For the Poisson ratio, the proportion of experimental observations exceeding the theoretical bounds is still smaller. It must be noted that a precise experimental evaluation of the Poisson ratio can be a challenge provided the much smaller range of variation when compared to the elastic moduli. Detailed results on the differences between the experimental values comprised in the training dataset and the theoretical bounds are shown in the Appendix A.

By comparing the values according to the hydration model, fewer points are outside the bounds when the TJ model is adopted for K or ν , which provides a more precise description of cement phases than the Powers model. These observations may suggest that the adoption of a precise description of cement paste phase assemblage is critical if theoretical bounds used to curate databases.

To conclude, experimental Young and shear moduli are in concordance with the bounds. For bulk modulus and Poisson ratio, only few points are in contradiction with the theoretical bounds. Depending on the trust given to the model in comparison with the experiments, it could be

decided to filter out the database of some experimental observations. However, here, for the proof of concept, all data are conserved to evaluate the ML performances without arbitration on the experimental results.

4.2 Prediction of elastic properties using ML and Micromechanics

ML and MB methods are evaluated to predict elastic properties of the samples contained in the training and test datasets.

4.2.1 Reproducing the training dataset observations

ML predictions. Knowing the w/c , DOH, and percentage fractions of clinker and gypsum in cement the four elastic properties are estimated by ML approaches. The validation procedure for one of the validation stages is illustrated in Figure 2.

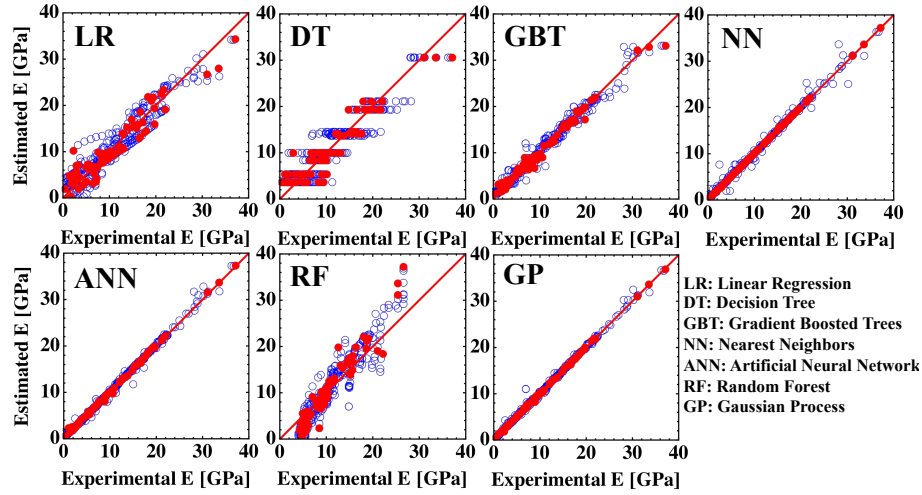


Figure 2: Illustration of the k -fold validation method (with $k = 5$): Predicted Young modulus E plotted against the experimental E at one validation stage out of 5 for the various ML methods tested: 292 values of the 4 training folds are depicted by empty blue dots, full red symbols depict the 73 elements used for validation.

The accuracy of predictions of elastic constants of cement pastes is compared for the various ML methods tested (Figure 3). The comparison serves to analyze the consistency and compatibility of the method regarding the database on which they are trained. The qualitative analysis suggests that the prediction of the Poisson ratio is less accurate when compared to predictions of the elastic moduli. Visually, NN, ANN, and GP perform better in predictions.

Errors are quantified using the root mean square error (RMSE):

$$\text{RMSE}(x) = \sqrt{\frac{\sum_{i=1}^n (x_i^{\text{pred}} - x_i^{\text{exp}})^2}{n}} \quad (9)$$

and mean relative error (MRE):

$$\text{MRE}(x) = \frac{1}{n} \sum_{i=1}^n \frac{|x_i^{\text{pred}} - x_i^{\text{exp}}|}{x_i^{\text{exp}}} \quad (10)$$

computed as a function of each prediction x_i^{pred} and experimental x_i^{exp} output averaged over the n observations i covering the whole training set obtained from the validation on all k -folds for the elastic constants. Tables 4 and 5 shows the RMSE and MRE, respectively, obtained for each ML method prediction. ANN, GP and NN yield the best accuracy in terms of RMSE and ME for the elastic constants.

MB estimations. Knowing the w/c , DOH, and percentage fractions of clinker and gypsum in cement the four elastic characteristics are also predicted by MB methods. Performances

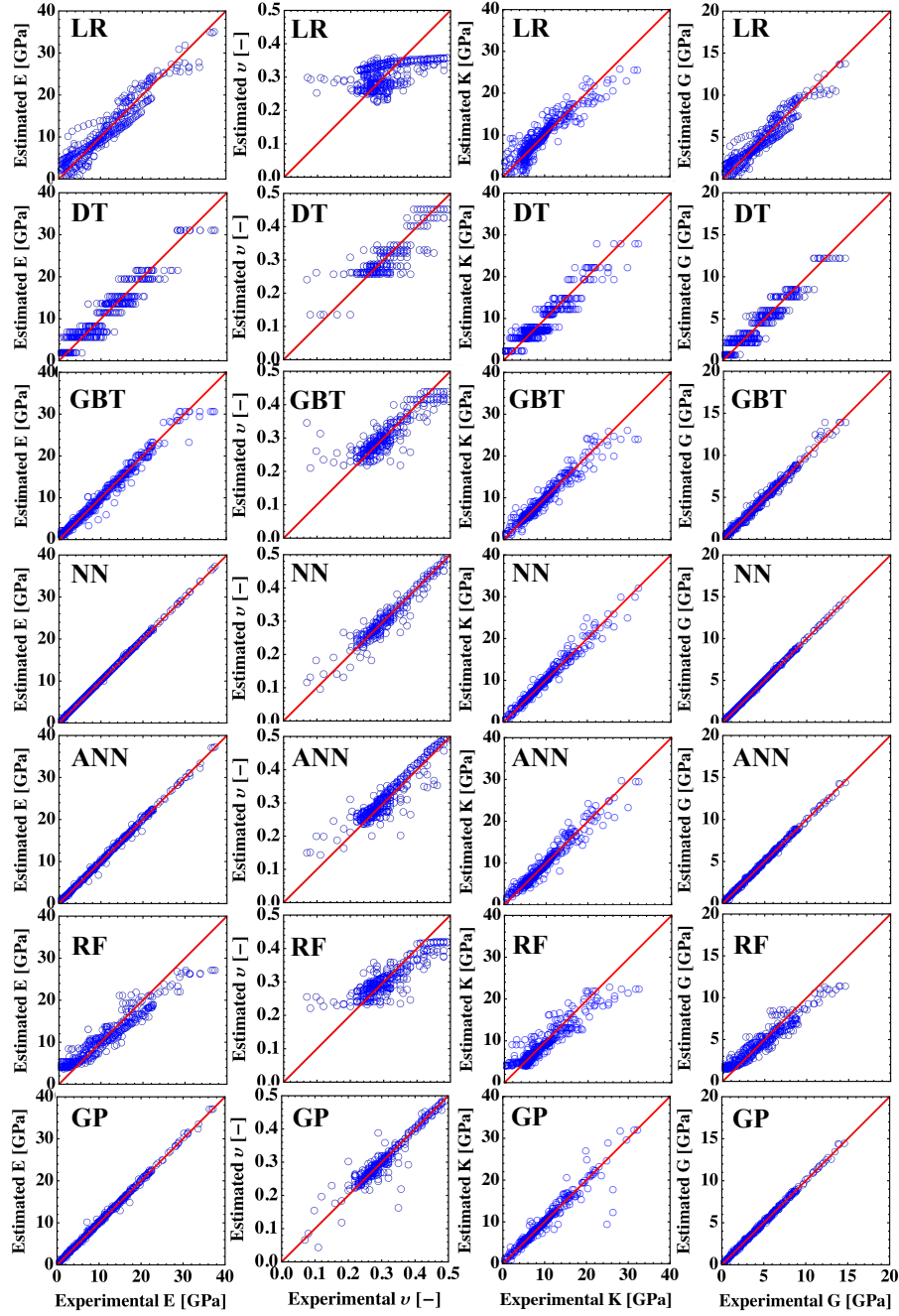


Figure 3: ML performances for elastic prediction of the training set: predicted values based on the various ML methods tested plotted against the experimental elastic constants from the training dataset (Young modulus E , Poisson ratio ν , bulk K and shear G moduli).

are shown in Figure 4. It can be noted that performances vary with the degree of hydration. The homogenization yields predictions of the elastic constants that are, in most cases, better when only the observations in the training dataset with $\text{DOH} \leq 0.7$ (i.e. associated with late ages) are accounted for (this effect can be more pronounced when MT estimates are used). The accuracy of MB estimations is quantified in Tables 6 and 7 using RMSE and MRE, respectively. These parameters were measured for the entire data set and also for the values in the training dataset with $\text{DOH} \geq 0.7$. When both error estimates are taken into consideration, the best MB estimated are given by SCPow And SC1s schemes for both cases when only spherical or ellipsoidal inclusions are considered. These four cases are used for comparison with ML methods. KHP model yields results closer to the ones obtained with Powers model.

As expected, when MB and ML methods predictions are confronted with the training dataset, ML methods display in general better accuracy than MB methods, with the less accurate ML

Table 4: Root mean square error (RMSE) of the elastic constants obtained from K-fold cross-validation technique based on 5-fold or 10-fold. The most accurate values are marked in bold.

Methods	RMSE(E) [GPa]		RMSE(ν) [-]		RMSE(K) [GPa]		RMSE(G) [GPa]	
	5-fold	10-fold	5-fold	10-fold	5-fold	10-fold	5-fold	10-fold
LR	2.4	2.4	0.061	0.062	2.2	2.2	1.0	1.0
DT	2.8	2.9	0.042	0.042	2.3	2.3	1.0	1.1
GBT	1.3	1.5	0.037	0.037	1.8	1.8	0.5	0.5
NN	1.1	1.1	0.034	0.032	1.6	1.5	0.4	0.4
ANN	0.6	0.6	0.034	0.033	1.5	1.4	0.3	0.2
RF	3.0	3.0	0.044	0.043	2.6	2.6	1.2	1.2
GP	0.7	0.8	0.032	0.031	2.0	2.9	0.3	0.2

Table 5: Mean relative error (MRE) of the elastic constants obtained from K-fold cross-validation technique based on 5-fold or 10-fold. The most accurate values are marked in bold.

Methods	MRE(E) [-]		MRE(ν) [-]		MRE(K) [-]		MRE(G) [-]	
	5-fold	10-fold	5-fold	10-fold	5-fold	10-fold	5-fold	10-fold
LR	0.50	0.50	0.18	0.18	0.34	0.34	0.54	0.56
DT	0.53	0.60	0.12	0.11	0.39	0.35	0.55	0.63
GBT	0.19	0.23	0.09	0.09	0.19	0.18	0.19	0.20
NN	0.17	0.16	0.07	0.07	0.13	0.13	0.18	0.17
ANN	0.08	0.10	0.08	0.07	0.16	0.15	0.13	0.11
RF	0.69	0.74	0.12	0.11	0.47	0.46	0.80	0.78
GP	0.10	0.09	0.06	0.06	0.17	0.22	0.12	0.10

methods (RF, DT) yielding predictions with a similar accuracy of the best MB estimations. The best RMSE and MRE for Young moduli predictions was with ANN and the accuracy was roughly 4-fold the accuracy of the best micromechanics estimation. Note however that the comparison according to the training dataset, of course, favors ML methods since these are trained specifically for them, whereas MB methods have not any *a priori* information on this correlation except the underlying theoretical model assumption.

4.2.2 Reproducing the test dataset observations

Predictions of elastic constants *for a test dataset* by ML methods and homogenization methods are also evaluated. As detailed in the Appendix A.1, the test dataset is composed of 58 observations including both static and dynamics measurements of elastic properties from various authors (Constantinides et al. 2004; Šavija et al. 2020; Chamrova 2010; Maruyama et al. 2014; Tamtsia et al. 2004; Haecker et al. 2005). The performance of MB and ML methods are analyzed in details for selected cases in the sequel.

Comparisons with data from Constantinides and Ulm (Constantinides et al. 2004) are given in Figure 5. For that case, ML methods (in particular, LR, ANN and DT) perform better in predicting experimental data than MB techniques. ML and MB yield comparable results when used to predict Haecker et al. (Haecker et al. 2005) (Figure 6) experimental data; the exceptions are the cases of homogenization methods using the Powers model. A similar result is obtained with Tamtsia et al. (Tamtsia et al. 2004) as visualized in Figure 7. In this case, homogenization with Tennis and Jennings hydration model performs quite well. On this test dataset, we can conclude that ML and MB have similar accuracy, none of them performs significantly better than the other.

The total RMSE and MRE associated with the test dataset for ML methods are shown in Figure 9 (values referring to the original training dataset l_{orig}) and will be discussed in the next section in comparison with prediction using extended training datasets.

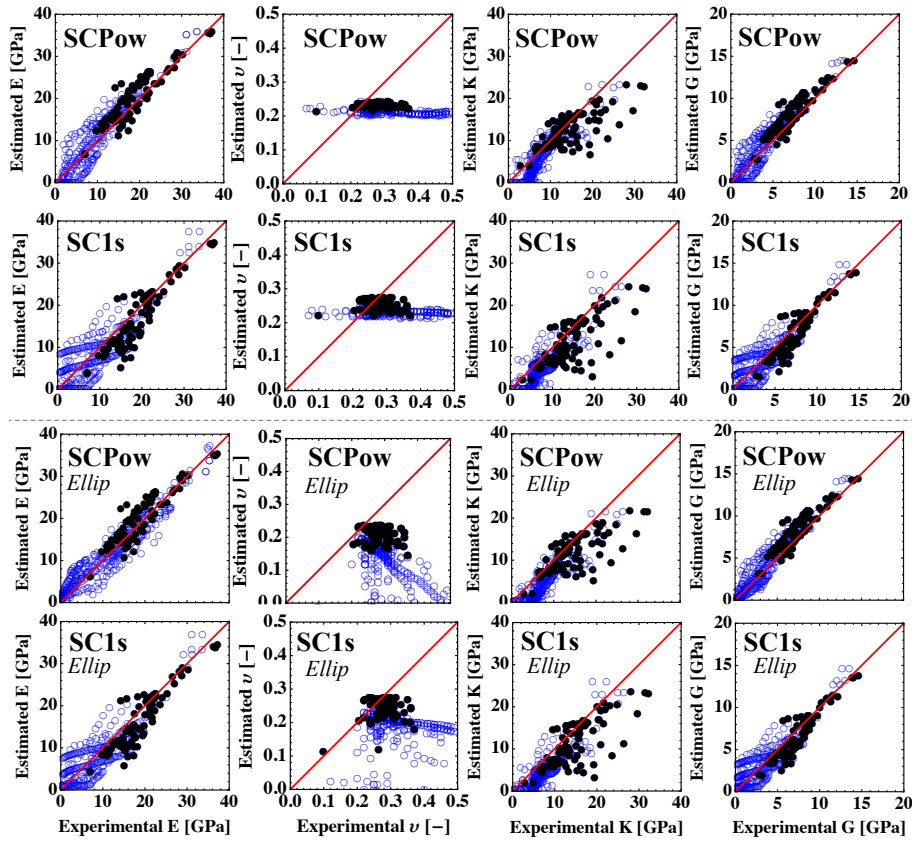


Figure 4: MB performances for elastic prediction of the training set (only the representations with the best performance are shown here): experimental elastic constants from the training dataset (Young modulus E , Poisson ratio ν , bulk K and shear G moduli) plotted against the estimated values using various homogenization methods. Influence of the degree of hydration on the performances: the empty blue circles correspond to $\text{DOH} \leq 0.7$, the solid black circles correspond to $\text{DOH} \geq 0.7$, i.e. predictions at late ages.

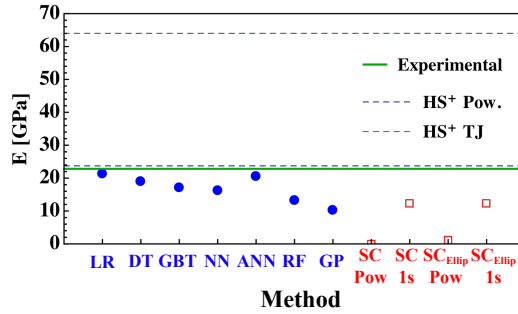


Figure 5: Prediction of the Young modulus E to reproduce the experimental observation by Constantinides and Ulm (Constantinides et al. 2004) (green solid line): results from ML (full blue dots) and MB (red empty dots) methods. Hashin-Shtrikman upper bounds using Powers and TJ model are shown for reference.

4.3 Missing data and extended database

ML and MB approaches can be used competitively, but ML approach can also be employed to guide experiments and obtain an optimized information over the whole parametric space and MB can be exploited to generate supplementary observations. Thus, direct experimental observations can be evaluated by ML and combined with synthetic observations obtained from MB schemes. Various methods are proposed in the literature to optimize the plan of experiments (Fuhg et al. 2020). Here, we adopt a distance-based approach to identify the zones in which a fewer number of experiments have been made. Then, we deploy a k -fold cross validation strategy to exploit domains in which metamodel interpolations are sufficient accurate and the domain in which

Table 6: Root mean square error (RMSE) of elastic constants as a measure of the accuracy of the homogenization estimations. The values in parenthesis correspond to estimations for elements in the dataset with $\text{DOH} \geq 0.7$ only. The most accurate values are marked in bold.

Methods	RMSE(E) [GPa]	RMSE(ν) [-]	RMSE(K) [GPa]	RMSE(G) [GPa]
MT Powers (MTPow)	9.7 (4.0)	0.10 (0.06)	4.5 (4.1)	4.1 (1.8)
SC Powers (SCPow)	2.8 (2.8)	0.11 (0.06)	3.4 (4.6)	1.2 (1.2)
MT KHP (MTKHP)	10.2 (7.2)	0.16 (0.06)	5.4 (5.8)	4.1 (2.9)
SC KHP (SCKHP)	8.1 (7.1)	0.10 (0.05)	6.4 (6.5)	3.2 (2.8)
MT 1-scale (MT1s)	6.9 (2.4)	0.09 (0.04)	3.6 (4.5)	2.8 (0.8)
SC 1-scale (SC1s)	3.8 (3.7)	0.10 (0.05)	3.8 (5.4)	1.5 (1.3)
MT 2-scales (MT2s)	6.9 (2.5)	0.09 (0.04)	3.6 (4.5)	2.8 (0.8)
SC 2-scales (SC2s)	7.6 (9.8)	0.08 (0.05)	6.7 (8.7)	2.9 (3.8)
MT _{Ellip} Powers (MTPow Ellip)	9.3 (3.9)	0.13 (0.06)	3.7 (4.4)	4.1 (1.8)
SC _{Ellip} Powers (SCPow Ellip)	2.5 (2.7)	0.30 (0.09)	4.2 (5.2)	1.1 (1.2)
MT _{Ellip} KHP (MTKHP Ellip)	12.7 (5.18)	0.08 (0.04)	9.0 (6.4)	5.1 (2.2)
SC _{Ellip} KHP (SCKHP Ellip)	3.2 (3.4)	0.30 (0.04)	5.0 (5.7)	1.6 (1.5)
MT _{Ellip} 1-scale (MT1sEllip)	6.6 (2.5)	0.11 (0.04)	3.2 (4.5)	2.8 (0.8)
SC _{Ellip} 1-scale (SC1s Ellip)	3.7 (3.6)	0.16 (0.06)	4.1 (5.5)	1.4 (1.3)
MT _{Ellip} 2-scales (MT2s Ellip)	6.6 (2.5)	0.11 (0.05)	3.2 (4.7)	2.8 (0.8)
SC _{Ellip} 2-scales (SC2s Ellip)	7.3 (9.1)	0.72 (0.44)	6.7 (8.2)	2.8 (3.4)

Table 7: Mean relative error (MRE) of elastic constants as a measure of the accuracy of the homogenization estimations. The values in parenthesis correspond to estimations for elements in the dataset with $\text{DOH} \geq 0.7$ only. The most accurate values are marked in bold.

Methods	MRE(E) [-]	MRE(ν) [-]	MRE(K) [-]	MRE(G) [-]
MT Powers (MTPow)	0.46 (0.23)	0.24 (0.18)	1.1 (0.22)	3.5 (0.27)
SC Powers (SCPow)	0.38 (0.14)	0.27 (0.19)	0.19 (0.38)	0.41 (0.15)
MT KHP (MTKHP)	3.6 (0.38)	0.52 (0.20)	0.94 (0.28)	2.5 (0.39)
SC KHP (SCKHP)	1.8 (0.38)	0.29 (0.14)	0.73 (0.33)	0.39 (0.77)
MT 1-scale (MT1s)	2.3 (0.11)	0.21 (0.12)	0.87 (0.21)	2.6 (0.10)
SC 1-scale (SC1s)	0.68 (0.18)	0.22 (0.13)	0.40 (0.27)	0.76 (0.17)
MT 2-scales (MT2s)	2.2 (0.10)	0.20 (0.12)	0.85 (0.20)	2.6 (0.09)
SC 2-scales (SC2s)	0.75 (0.53)	0.20 (0.15)	0.72 (0.56)	0.77 (0.53)
MT _{Ellip} Powers (MTPow)	0.45 (0.22)	0.32 (0.20)	0.87 (0.22)	3.5 (0.27)
SC _{Ellip} Powers (SCPow)	0.37 (0.13)	0.75 (0.27)	0.47 (0.22)	0.40 (0.16)
MT _{Ellip} KHP (MTKHP Ellip)	0.66 (0.19)	0.18 (0.12)	1.6 (0.30)	4.1 (0.24)
SC _{Ellip} KHP (SCKHP Ellip)	2.0 (0.13)	0.39 (0.12)	0.52 (0.26)	0.41 (0.18)
MT _{Ellip} 1-scale (MT1s)	2.2 (0.11)	0.26 (0.12)	0.68 (0.20)	2.6 (0.10)
SC _{Ellip} 1-scale (SC1s)	0.61 (0.17)	0.37 (0.14)	0.42 (0.28)	0.71 (0.16)
MT _{Ellip} 2-scales (MT2s)	2.2 (0.11)	0.26 (0.13)	0.67 (0.21)	2.5 (0.09)
SC _{Ellip} 2-scales (SC2s)	0.70 (0.49)	0.73 (1.55)	0.73 (0.53)	0.68 (0.48)

information in lacking.

4.3.1 Distance-based approach to identify the domains with missing data

To identify the domains with missing data, we adopt a simple strategy based on the distance of data in a given dimension of input space. In each dimension of the input space, we order the components of the observations in ascending order as shown in Figure 8(top). The normalized

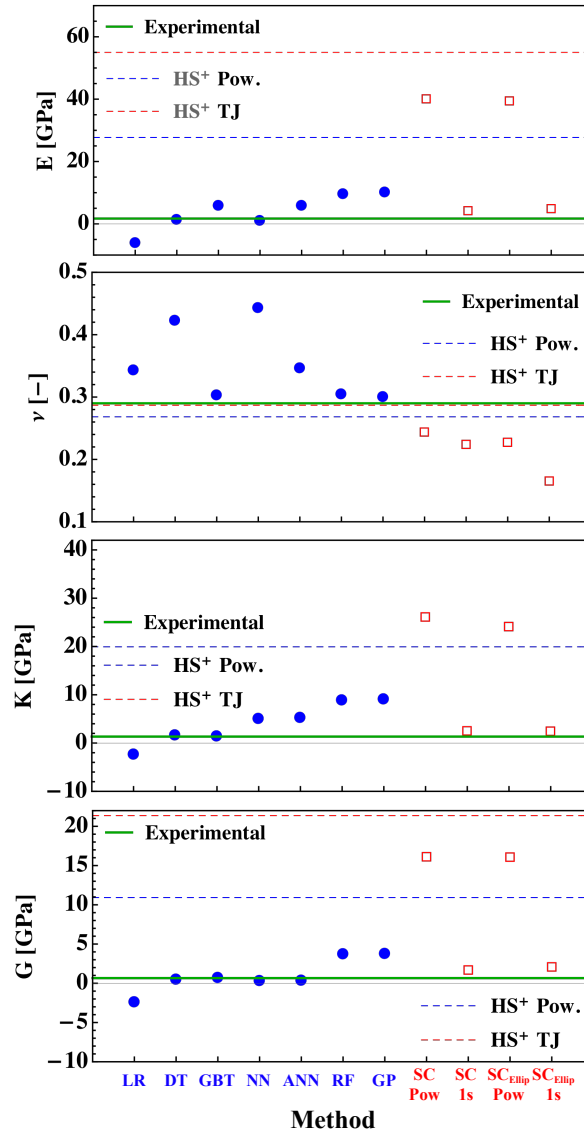


Figure 6: Prediction of the elastic constants (Young modulus E , Poisson ratio ν , shear G and bulk K moduli) to reproduce the experimental observations obtained by Haecker et al. (Haecker et al. 2005) (green solid line); results from ML (full blue dots) and MB (red empty squares) methods. Hashin-Shtrikman upper bounds using Powers and TJ model are shown for reference.

difference between two component $x_{(i'+1)}$ and $x_{i'}$: $\Delta O_N = \frac{1}{\sum_i x_i} (x_{(i'+1)} - x_{i'})$ is related to the extent of the domain associated with missing data, where i' denotes the ordered position of an observation. Figure 8 (bottom) shows $\Delta O_N(i')$ for all input components considered here. A large $\Delta O_N(i')$ indicates that two ordered observations $x_{(i'+1)}$ and $x_{i'}$ are relatively far from each other and that the interval $]x_{(i')}, x_{(i'+1)}[$ is a zone in which data is missing. To identify the most relevant zones in which data is missing according to this approach, we adopt the following criterion: a new observation $x_i^* = \frac{1}{2}(x_{i'} + x_{(i'+1)})$ is to be generated whenever $\Delta O_N \geq c_o$, where c_o is an arbitrary cut-off. We adopt (gray dashed lines in Figure 8 (bottom)) $c_o=0.0001$ for the w/c and degree of hydration, and $c_o=0.0005$ for the clinker minerals and gypsum mass fractions. With this approach, the selected x_i^* per input are:

- for $j = DOH$ [-]: 0.70, 0.95
- for $j = w/c$ [-]: 0.28, 0.37, 0.72
- for $j = m_{C_3S}$ [%]: 33.05, 52.405, 85.25
- for $j = m_{C_2S}$ [%]: 5.40, 11.24, 16.70, 20.70, 31.75, 49.95
- for $j = m_{C_3A}$ [%]: 1.10, 9.95
- for $j = m_{C_4AF}$ [%]: 4.40, 11.25

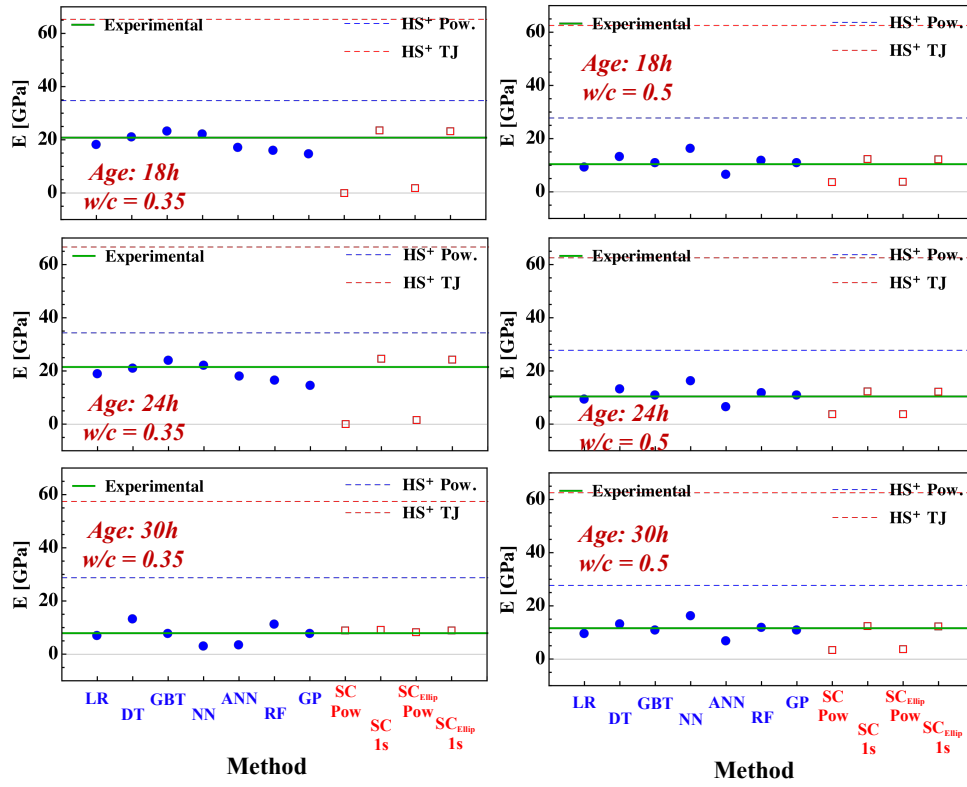


Figure 7: Prediction of the Young modulus E for $w/c = 0.35$ and $w/c = 0.50$ at different ages to reproduce the experimental observations by Tamtsia et al. (Tamtsia et al. 2004) (green solid line): results from ML (full blue dots) and MB (empty red squares) methods. Hashin-Shtrikman upper bounds using Powers and TJ model are shown for reference.

- for $j = m_{Gypsum} [\%]$: 1.05, 2.00, 2.90, 3.65, 4.69, 6.09

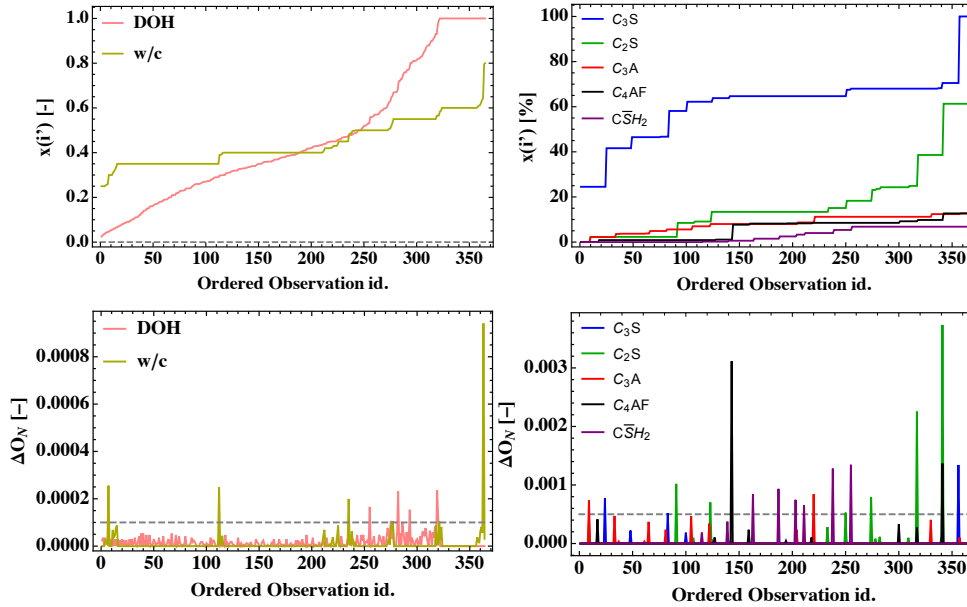


Figure 8: Strategy to identify the zones with missing data. At the top, ordered input vector $x(i')$ for each input component \in Age, w/c , m_{C_3S} , m_{C_2S} , m_{C_3A} , m_{C_4AF} , m_{Gypsum} as a function of the ordered observation index i' . At the bottom, normalized difference between two component $x_{(i'+1)}$ and $x_{i'}$: $\Delta O_N = \frac{1}{\sum_i x_i} (x_{(i'+1)} - x_{i'})$ as a function of index i' indicating ascending order per component. The gray dashed lines depict the limit criterion adopted to identify the most relevant domains with data missing.

To generate the new data in these zones, we defined three new datasets:

- The minimum dataset l_{min} covering all x_i^* for all input vector O_m identified by the strategy

above. The minimum number of observations to be generated covering all these values is 6.

- A dataset l_{exis}^{1P} with one new observation by O_{ij}^* per input identified, with each one of the 24 O_{ij}^* values associated with an *already existing* (randomly sampled) set of inputs.
- A dataset l_{self}^{1P} also with 24 observations with each observation being a random combination of the O_{ij}^* identified by the strategy above.

To generate the new data on elastic constants, we adopt the MB method SC1s_{Ellip}, which yield one of the best performances of MB methods, as discussed in Section 4.2.1.

In Figure 9, we compare the performance of ML methods trained on the original and extended datasets on the estimation of the Young modulus of the test dataset (RMSE and ME at the left), and training dataset *via* a cross-validation approach (k=5 folds) on the training dataset (RMSE and ME at the right).

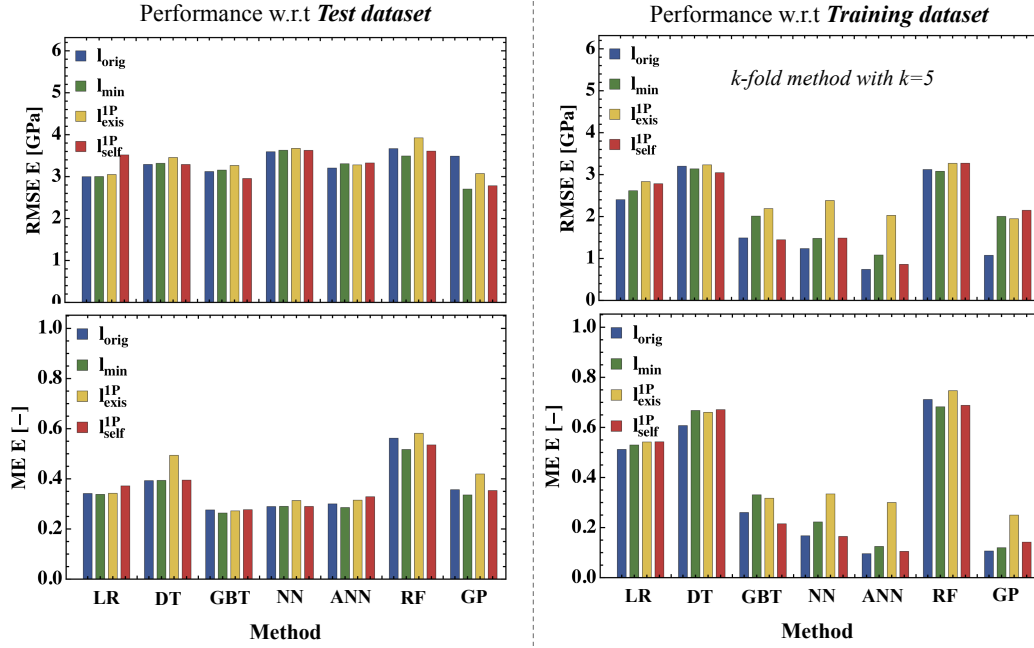


Figure 9: RMSE and ME computed for the test dataset (left) and using a cross-validation approach (k=5 folds) on the training dataset (right) for various ML methods using only the original training dataset l_{orig} , or extended training datasets l_{min} , l_{exis}^{1P} , or l_{self}^{1P} . New data is generated with SC1s_{Ellip}.

Regarding the performance when the test dataset is considered, the use of the extended training dataset l_{exis}^{1P} improves the accuracy of ANN predictions, and the use of l_{exis}^{2P} improves the accuracy of GP predictions. Regarding the performance when the training dataset is considered, the use of extended training datasets generally increase the error as computed by the k-fold methods except for DT, GBT and RF. When the extended training datasets lead to worse accuracies, it must be noted that the increase in RMSE and ME is not too large (except in some of the l_{exis}^{1P} cases in which the error is duplicated when the training dataset is considered). This observation suggests that using MB methods to generate missing data does not impair significantly the precision of predictions.

These results show that MB methods can be used to generate new data to complete databases for establishing composition-property correlations in cement-based materials leading in some cases (ANN and GP, the best performances in prediction) to an improvement in the prediction accuracy regarding the test dataset.

5 Conclusions

In this article, Machine Learning and Micromechanics-based methods were deployed to establish correlations between the composition and the elastic property of OPC pastes. In the exploration of the methods, we identified opportunities for using them as promising allies. ML arises as a proficient tool to exploit variety of results from different authors with a set of input characteristics, and identify significant lack of knowledge. Micromechanics is an opportunity to judge experimental

results, provides bounds to check ML predictions, and furnish supplementary/complementary data for ML to be trained on. The main conclusions of this study are as follows:

- *On the methods to link composition and elastic properties.* The accuracy of ML and MB predictions are comparable for predicting elastic properties of the *test dataset*. When MB and ML are confronted with the *training dataset*, ML gives better accuracy than MB methods, with the less accurate ML methods (RF, DT) yielding predictions with a similar accuracy of the best MB estimations (the comparison according to the training dataset, of course, favors ML methods since these are trained specifically for them, whereas MB methods have not any *a priori* information on this correlation except the underlying theoretical model assumption.). It must be noted that the test dataset used in this work use both dynamic and static measurements while the ML methods were training only in dynamic experimental data. Both ML and analytical micromechanics computations performed here are not computer-intensive, especially when compared to often fastidious numerical homogenization approaches. This aspect is a clear advantage of ML and analytical MB methods, notably when a larger exploration of the compositional design space is desired.
- *Data-driven estimates and importance of reliable databases.* Even with a relatively small dataset, ML methods have proven to be reliable and robust in the prediction of elastic properties of cement pastes from their composition for test dataset. Thus, the effort to build and enlarge the databases on cement composition and properties, for instance including static measurements in the training data set or even using the same strategy for other properties than elastic properties, may benefit cement and concrete research providing a reliable tool to tailor the composition of the material for a target property or performance specification.
- *Providing missing data.* Analytical micromechanics methods appear proficient to complete the database for input values which have not been explored by experimental campaigns. Indeed, the accuracy of ML and MB being comparable corroborates that MB methods can be used to provide missing data in the databases of cement-based materials despite their well-known variability, the significant lack of knowledge being robustly identified from the ML approaches. This observation adds to the accumulating evidence showing that MB approaches are a powerful tool to estimate the property from the composition based on a few fundamental component data set and assumptions on cement hydration and micro-structure models. Besides, providing virtual estimations for concealing missing experimental data, they can also serve to crosscheck uncertain or suspicious observations.

The strategy outlined in this study combining MB and ML methods to explore the space of formulation design linking it to the effective properties of the materials can be extended to other properties in cement and concrete science. This approaches arises an interesting and not costly tool to feed mechanical simulations taking into account the local variability of material properties based on physical and micro-scale properties, even for large simulations. It has been exploited for domains with missing data, it could also be enriched with analysis of subdomains with non-trustable data.

6 Bibliography

- Achour, M., F. Bignonnet, J.-F. Barthélémy, E. Rozière, and O. Amiri (Feb. 2020). "Multi-scale modeling of the chloride diffusivity and the elasticity of Portland cement paste". *Construction and Building Materials* 234, p. 117124
- Acker, P. (2001). "Micromechanical analysis of creep and shrinkage mechanisms." *Creep, shrinkage and durability mechanics of concrete and other quasi-brittle materials*.
- Agrawal, A. and A. Choudhary (Apr. 2016). "Perspective: Materials informatics and big data: Realization of the "fourth paradigm" of science in materials science". *APL Materials* 4.5, p. 053208
- Aller, L. H., I. Appenzeller, B. Baschek, K. Butler, C. De Loore, H. W. Duerbeck, M. F. El Eid, H. H. Fink, T. Herczeg, T. Richtler, H. Schneider, M. Scholz, W. Seggewiss, W. C. Seitter, J. Trümper, P. Ulmenschneider, R. Wehrse, and V. Weidemann (1996). *Landolt-Börnstein: Numerical Data and Functional Relationships in Science and Technology - New Series " Gruppe/Group 6*

- Astronomy and Astrophysics " Volume 3 Voigt: Astronomy and Astrophysics. Extension and Supplement to Volume 2 " Stars and Star Clusters.* Iborlbor
- Bary, B. and S. Béjaoui (Feb. 2006). "Assessment of diffusive and mechanical properties of hardened cement pastes using a multi-coated sphere assemblage model". *Cement and Concrete Research* 36.2, pp. 245–258
- Behnood, A., J. Olek, and M. A. Glinicki (Sept. 2015). "Predicting modulus elasticity of recycled aggregate concrete using M5 model tree algorithm". *Construction and Building Materials* 94, pp. 137–147
- Ben Chaabene, W., M. Flah, and M. L. Nehdi (Nov. 2020). "Machine learning prediction of mechanical properties of concrete: Critical review". *Construction and Building Materials* 260, p. 119889
- Bengio, Y. and Y. Grandvalet (2004). "No Unbiased Estimator of the Variance of K-Fold Cross-Validation". *Journal of Machine Learning Research* 5, pp. 1089–1105
- Boumiz, A., C. Vernet, and F. C. Tenoudji (Apr. 1996). "Mechanical properties of cement pastes and mortars at early ages: Evolution with time and degree of hydration". *Advanced Cement Based Materials* 3.3, pp. 94–106
- Breiman, L. (Aug. 1996). "Bagging predictors". *Machine Learning* 24.2, pp. 123–140
- Breiman, L., J. H. Friedman, R. A. Olshen, and C. J. Stone (1984). "Classification and regression trees. Belmont, CA: Wadsworth". *International Group* 432, pp. 151–166
- Bullard, J. W., E. J. Garboczi, P. E. Stutzman, P. Feng, A. S. Brand, L. Perry, J. Hagedorn, W. Griffin, and J. E. Terrill (July 2017). "Measurement and modeling needs for microstructure and reactivity of next-generation concrete binders". *Cement and Concrete Composites*
- Chamrova, R. (2010). "Modelling and measurement of elastic properties of hydrating cement paste". PhD Thesis. EPFL
- Constantinides, G. and F.-J. Ulm (Jan. 2004). "The effect of two types of C-S-H on the elasticity of cement-based materials: Results from nanoindentation and micromechanical modeling". *Cement and Concrete Research* 34.1, pp. 67–80
- Duan, Z. H., S. C. Kou, and C. S. Poon (Mar. 2013). "Prediction of compressive strength of recycled aggregate concrete using artificial neural networks". *Construction and Building Materials* 40, pp. 1200–1206
- Fuhg, J. N., A. Fau, and U. Nackenhorst (2020). "State-of-the-Art and Comparative Review of Adaptive Sampling Methods for Kriging". *Archives of Computational Methods in Engineering*, pp. 1–59. DOI: [10.1007/s11831-020-09474-6](https://doi.org/10.1007/s11831-020-09474-6). URL: <https://link.springer.com/article/10.1007/s11831-020-09474-6>
- Ghabezloo, S., J. Sulem, and J. Saint-Marc (Jan. 2009). "The effect of undrained heating on a fluid-saturated hardened cement paste". *Cement and Concrete Research* 39.1, pp. 54–64
- Golafshani, E. M. and A. Behnood (Mar. 2018). "Automatic regression methods for formulation of elastic modulus of recycled aggregate concrete". *Applied Soft Computing* 64, pp. 377–400
- Guihard, V., F. Taillade, J.-P. Balayssac, B. Steck, and J. Sanahuja (July 2019). "Permittivity measurement of cementitious materials and constituents with an open-ended coaxial probe: combination of experimental data, numerical modelling and a capacitive model". *RILEM Technical Letters* 4, pp. 39–48
- Haecker, C. .-, E. J. Garboczi, J. W. Bullard, R. B. Bohn, Z. Sun, S. P. Shah, and T. Voigt (Oct. 2005). "Modeling the linear elastic properties of Portland cement paste". *Cement and Concrete Research* 35.10, pp. 1948–1960
- Hammond, C. (1997). *Handbook of chemistry and physics*. Lide DR, CRC Press, Boca Raton, FL
- Hansen, T. C. (Nov. 1986). "Physical structure of hardened cement paste. A classical approach". en. *Materials and Structures* 19.6, pp. 423–436. DOI: [10.1007/BF02472146](https://doi.org/10.1007/BF02472146). (Visited on 11/15/2020)
- Hashin, Z. and S. Shtrikman (Mar. 1963). "A variational approach to the theory of the elastic behaviour of multiphase materials". *Journal of the Mechanics and Physics of Solids* 11.2, pp. 127–140. DOI: [10.1016/0022-5096\(63\)90060-7](https://doi.org/10.1016/0022-5096(63)90060-7). (Visited on 10/29/2014)
- Helmuth, R. A. and D. H. Turk (1966). "Elastic moduli of hardened Portland cement and tricalcium silicate pastes: effect of porosity". *Special Report - Highway Research Board* 90, p. 135
- Hlobil, M. (June 2020). "Distribution of hydration products in the microstructure of cement pastes". *Acta Polytechnica CTU Proceedings* 27.0, pp. 84–89

- Honorio, T., B. Bary, and F. Benboudjema (July 2016). "Multiscale estimation of ageing viscoelastic properties of cement-based materials: A combined analytical and numerical approach to estimate the behaviour at early age". *Cement and Concrete Research* 85, pp. 137–155
- Honorio, T., B. Bary, and F. Benboudjema (Mar. 2018). "Thermal properties of cement-based materials: Multiscale estimations at early-age". *Cement and Concrete Composites* 87, pp. 205–219
- Honorio, T., B. Bary, F. Benboudjema, and S. Poyet (May 2016). "Modeling hydration kinetics based on boundary nucleation and space-filling growth in a fixed confined zone". *Cement and Concrete Research* 83, pp. 31–44
- Honorio, T., T. Bore, F. Benboudjema, E. Vourc'h, and M. Ferhat (Mar. 2020). "Dielectric properties of the pore solution in cement-based materials". *Journal of Molecular Liquids* 302, p. 112548
- Honorio, T., L. Brochard, and B. Bary (Aug. 2018). "Statistical variability of mechanical fields in thermo-poro-elasticity: Multiscale analytical estimations applied to cement-based materials at early-age". *Cement and Concrete Research* 110, pp. 24–41
- Honorio, T., H. Carasek, and O. Cascudo (June 2020). "Electrical properties of cement-based materials: Multiscale modeling and quantification of the variability". *Construction and Building Materials* 245, p. 118461
- Honorio, T., P. Guerra, and A. Bourdot (Sept. 2020). "Molecular simulation of the structure and elastic properties of ettringite and monosulfoaluminate". *Cement and Concrete Research*, p. 106126
- Inc., W. R. (n.d.). *Mathematica, Version 13.0.0*. Champaign, IL, 2021. URL: <https://www.wolfram.com/mathematica>
- Kasperkiewicz, J., J. Racz, and A. Dubrawski (Oct. 1995). "HPC Strength Prediction Using Artificial Neural Network". *Journal of Computing in Civil Engineering* 9.4, pp. 279–284
- Ke, G., Q. Meng, T. Finley, T. Wang, W. Chen, W. Ma, Q. Ye, and T.-Y. Liu (2017). "LightGBM: A Highly Efficient Gradient Boosting Decision Tree". *Advances in Neural Information Processing Systems 30 (NIPS 2017)*, p. 9
- Königsberger, M., C. Hellmich, and B. Pichler (Oct. 2016). "Densification of C-S-H is mainly driven by available precipitation space, as quantified through an analytical cement hydration model based on NMR data". *Cement and Concrete Research* 88, pp. 170–183. DOI: [10.1016/j.cemconres.2016.04.006](https://doi.org/10.1016/j.cemconres.2016.04.006). URL: <http://www.sciencedirect.com/science/article/pii/S0008884616303374> (visited on 08/05/2016)
- Königsberger, M., T. Honório, J. Sanahuja, B. Delsaute, and B. L. A. Pichler (July 5, 2021). "Homogenization of nonaging basic creep of cementitious materials: A multiscale modeling benchmark". *Construction and Building Materials* 290, p. 123144. DOI: [10.1016/j.conbuildmat.2021.123144](https://doi.org/10.1016/j.conbuildmat.2021.123144). (Visited on 04/30/2021)
- Lura, P., O. M. Jensen, and K. van Breugel (Feb. 2003). "Autogenous shrinkage in high-performance cement paste: An evaluation of basic mechanisms". *Cement and Concrete Research* 33.2, pp. 223–232
- Manzano, H. (2009). "Atomistic Simulation studies of the Cement Paste Components". PhD thesis. Universidad del País Vasco
- Maruyama, I. and G. Igarashi (2014). "Cement Reaction and Resultant Physical Properties of Cement Paste". *Journal of Advanced Concrete Technology* 12.6, pp. 200–213
- Monteiro, P. J. M. and C. T. Chang (Dec. 1995). "The elastic moduli of calcium hydroxide". *Cement and Concrete Research* 25.8, pp. 1605–1609
- Muller, A. C. A., K. L. Scrivener, A. M. Gajewicz, and P. J. McDonald (2012). "Densification of C–S–H Measured by ^1H NMR Relaxometry". *The Journal of Physical Chemistry C* 117.1, pp. 403–412. DOI: [10.1021/jp3102964](https://doi.org/10.1021/jp3102964). URL: <http://dx.doi.org/10.1021/jp3102964> (visited on 12/19/2014)
- Mura, T. (1987). *Micromechanics of defects in solids*. 2nd ed. Mechanics of elastic and inelastic solids. Editors: S. Nemat-Nasser and G. AE. Oravas. Martinus Nijhoff Publishers
- Nematzadeh, Z., R. Ibrahim, and A. Selamat (May 2015). "Comparative studies on breast cancer classifications with k-fold cross validations using machine learning techniques". *2015 10th Asian Control Conference (ASCC)*. DOI: [10.1109/ASCC.2015.7244654](https://doi.org/10.1109/ASCC.2015.7244654)

- Olson, G. B. (Aug. 1997). "Computational Design of Hierarchically Structured Materials". *Science* 277.5330, pp. 1237–1242
- Patel, R. A., Q. T. Phung, S. C. Seetharam, J. Perko, D. Jacques, N. Maes, G. De Schutter, G. Ye, and K. Van Breugel (Dec. 2016). "Diffusivity of saturated ordinary Portland cement-based materials: A critical review of experimental and analytical modelling approaches". *Cement and Concrete Research* 90, pp. 52–72
- Pichler, B. and C. Hellmich (May 2011). "Upscaling quasi-brittle strength of cement paste and mortar: A multi-scale engineering mechanics model". *Cement and Concrete Research* 41.5, pp. 467–476
- Pichler, B., C. Hellmich, J. Eberhardsteiner, J. Wasserbauer, P. Termkhajornkit, R. Barbarulo, and G. Chanvillard (Mar. 2013). "Effect of gel–space ratio and microstructure on strength of hydrating cementitious materials: An engineering micromechanics approach". *Cement and Concrete Research* 45, pp. 55–68
- Powers, T. C. (1960). "Physical properties of cement paste". *Proceedings of the Fourth International Symposium on Chemistry of Cement*. Washington, pp. 577–613
- Powers, T. C. and T. L. Brownyard (Sept. 1946). "Studies of the Physical Properties of Hardened Portland Cement Paste". *Journal Proceedings* 43.9, pp. 101–132
- Rodriguez, J. D., A. Perez, and J. A. Lozano (Mar. 2010). "Sensitivity Analysis of k-Fold Cross Validation in Prediction Error Estimation". *IEEE Transactions on Pattern Analysis and Machine Intelligence* 32.3. Conference Name: IEEE Transactions on Pattern Analysis and Machine Intelligence, pp. 569–575. DOI: [10.1109/TPAMI.2009.187](https://doi.org/10.1109/TPAMI.2009.187)
- Sanahuja, J., L. Dormieux, and G. Chanvillard (2007). "Modelling elasticity of a hydrating cement paste". *Cement and Concrete Research* 37.10, pp. 1427–1439
- Šavija, B., H. Zhang, and E. Schlangen (Apr. 2020). "Micromechanical testing and modelling of blast furnace slag cement pastes". en. *Construction and Building Materials* 239, p. 117841. DOI: [10.1016/j.conbuildmat.2019.117841](https://doi.org/10.1016/j.conbuildmat.2019.117841). URL: <http://www.sciencedirect.com/science/article/pii/S0950061819332945> (visited on 11/10/2020)
- Speziale, S., F. Jiang, Z. Mao, P. J. M. Monteiro, H.-R. Wenk, T. S. Duffy, and F. R. Schilling (July 2008). "Single-crystal elastic constants of natural ettringite". *Cement and Concrete Research* 38.7, pp. 885–889
- Sun, Z., E. J. Garboczi, and S. P. Shah (Jan. 2007). "Modeling the elastic properties of concrete composites: Experiment, differential effective medium theory, and numerical simulation". *Cement and Concrete Composites* 29.1, pp. 22–38
- Tamtsia, B. T., J. J. Beaudoin, and J. Marchand (July 2004). "The early age short-term creep of hardening cement paste: load-induced hydration effects". *Cement and Concrete Composites* 26.5, pp. 481–489
- Tennis, P. D. and H. M. Jennings (2000). "A model for two types of calcium silicate hydrate in the microstructure of Portland cement pastes". *Cement and Concrete Research* 30, pp. 855–863
- Termkhajornkit, P. and R. Barbarulo (2012). "Modeling the coupled effects of temperature and fineness of Portland cement on the hydration kinetics in cement paste". *Cement and Concrete Research* 42.3, pp. 526–538
- Torquato, S. (2002). *Random Heterogeneous Materials: Microstructure and Macroscopic Properties*. Springer Science & Business Media
- Ulm, F.-J., G. Constantinides, and F. H. Heukamp (Jan. 2004). "Is concrete a poromechanics materials?—A multiscale investigation of poroelastic properties". *Materials and Structures* 37.1, pp. 43–58
- Velez, K., S. Maximilien, D. Damidot, G. Fantozzi, and F. Sorrentino (Apr. 2001). "Determination by nanoindentation of elastic modulus and hardness of pure constituents of Portland cement clinker". *Cement and Concrete Research* 31.4, pp. 555–561
- Wang, X. and K. V. Subramaniam (Mar. 2011). "Ultrasonic monitoring of capillary porosity and elastic properties in hydrating cement paste". *Cement and Concrete Composites* 33.3, pp. 389–401
- Wangler, T., N. Roussel, F. P. Bos, T. A. M. Salet, and R. J. Flatt (Sept. 2019). "Digital Concrete: A Review". *Cement and Concrete Research* 123, p. 105780

- Williams, C. K. I. and C. E. Rasmussen (1996). “Gaussian Processes for Regression”. *Advances in neural information processing systems*, p. 7
- Wyrzykowski, M., J. Sanahuja, L. Charpin, M. Königsberger, C. Hellmich, B. Pichler, L. Valentini, T. Honório, V. Smilauer, K. Hajkova, G. Ye, P. Gao, C. Dunant, A. Hilaire, S. Bishnoi, and M. Azenha (Dec. 2017). “Numerical benchmark campaign of COST Action TU1404 – microstructural modelling”. *RILEM Technical Letters* 2, pp. 99–107
- Yan, K. and C. Shi (Aug. 2010). “Prediction of elastic modulus of normal and high strength concrete by support vector machine”. *Construction and Building Materials* 24.8, pp. 1479–1485
- Yeh, I.-C. and L.-C. Lien (Apr. 2009). “Knowledge discovery of concrete material using Genetic Operation Trees”. *Expert Systems with Applications* 36.3, Part 2, pp. 5807–5812
- Yeh, I. -. (Dec. 1998). “Modeling of strength of high-performance concrete using artificial neural networks”. *Cement and Concrete Research* 28.12, pp. 1797–1808
- Young, B. A., A. Hall, L. Pilon, P. Gupta, and G. Sant (Jan. 2019). “Can the compressive strength of concrete be estimated from knowledge of the mixture proportions?: New insights from statistical analysis and machine learning methods”. *Cement and Concrete Research* 115, pp. 379–388
- Zaoui, A. (Aug. 2002). “Continuum Micromechanics: Survey”. *Journal of Engineering Mechanics* 128.8, pp. 808–816. DOI: [10.1061/\(ASCE\)0733-9399\(2002\)128:8\(808\)](https://doi.org/10.1061/(ASCE)0733-9399(2002)128:8(808)). (Visited on 10/14/2020)
- Zimmerman, R. W. (1992). “Hashin-Shtrikman bounds on the poisson ratio of a composite material”. *Mechanics Research Communications* 19.6, pp. 563–569. DOI: [10.1016/0093-6413\(92\)90085-O](https://doi.org/10.1016/0093-6413(92)90085-O). (Visited on 03/19/2021)

7 Supplementary material

The database with observations on cement pastes with cement composition and elastic constants is available at https://github.com/tuliohf/cdi/blob/main/Dataset_Elastic_Constants OPC_Pastes

Competing interests

- The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.
- The authors declare that they have no competing interests.

Appendix A. Database collection and analysis

Experimental observations, as detailed below, are collected from various papers, and results are analyzed in comparison to some standard modelling results.

A.1 Database collection

The database collected comprises 365 observations with “full” information on which ML methods are trained, and 11 observations with “partial” information used for testing both ML and MB approaches.

A.1.1. Training dataset

Only data from dynamical measurements of elastic constants are considered in the training dataset.

Helmuth and Turk (1966) (Helmuth et al. 1966). The authors do not provide the w/c ratio chosen for experimental formulation but instead they measure the ratio w_t/c_i between the total water content w_t and the ignited weight c_i at late ages. As proposed by other authors (*Achour et al. 2020; Sanahuja et al. 2007*), it is possible to estimate the w/c ratio using

$$w_t/c_i = \begin{cases} w/c \frac{\kappa_w}{\kappa_h - 1}, & \text{if } w/c \leq \frac{\kappa_h - 1}{\rho_{Clinker}} \\ w/c + \frac{1 + \kappa_w - \kappa_h}{\rho_{Clinker}}, & \text{otherwise} \end{cases}$$

where $\rho_{Clinker} = 3.13$ is the density of clinker. The quantities $\kappa_h = 2.13$ and $\kappa_w = 1.31$ are the volumes of depleted water and formed hydrates, respectively, per one unit of clinker consumed by hydration processes.

Boumiz et al. (1996) ([Boumiz et al. 1996](#)). Only data on cement pastes were published. For some observations, elastic properties were provided without knowledge about age or degree of hydration. In these cases, missing key information has been approximated by local linear least-squares regression. Since the experimental data is pretty smooth, approximating the local behavior (in the range of a few observation points) by a linear fitting is a reasonable choice.

Haecker et al. (2005) ([Haecker et al. 2005](#)). The data on cements “H” and “D” were collected as presented by the authors.

Sun et al. (2007) ([Sun et al. 2007](#)). As for Boumiz et al. (1996), local linear least-squared regression was performed to obtain age and degree of hydration for some experimental observations. Modified Bogue formula was used to compute clinker mineral fractions. Besides elastic constants of cement pastes, the same study reports also results at mortar and concrete scales that can be used in future work.

Wang and Subramaniam (2011) ([Wang et al. 2011](#)). As for Boumiz et al. (1996), local linear least-squared regression was performed to obtain age and degree of hydration for some experimental observations. Modified Bogue formula was used to compute clinker mineral fractions.

Chamrova (2010) ([Chamrova 2010](#)). The data were collected as presented by the authors.

Maruyama and Igarashi (2014) ([Maruyama et al. 2014](#)). As for Boumiz et al. (1996), the age and degree of hydration were estimated for some observations by local linear least-squares regression.

A.1.2. Test dataset

Experimental observations from ([Tamtsia et al. 2004](#); [Constantinides et al. 2004](#); [Lura et al. 2003](#); [Šavija et al. 2020](#)) did not include information regarding either age, degree of hydration, or the pair of elastic constants necessary to characterize isotropic elastic behavior. Therefore, these samples could not be included in the training dataset, they form the core of the test dataset. The data on cement “L” from ([Haecker et al. 2005](#)) is also incorporated in the test dataset.

For the the test dataset, both static and dynamical measurements are considered indistinctly.

A.2 Database analysis

The experimental observations are compared with the theoretical bounds of elastic properties provided by analytical estimations for random heterogeneous media. Voigt-Reuss bounds and the Hashin-Shtrikman bounds, as introduced in Section 3.4, are considered. Figure 10 shows the differences for each observation included in the training set between the experimental values of the elastic properties and the upper Voigt and Hashin-Shtrikman bounds obtained from the knowledge of the phase fractions and the homogenization schemes. The positive values correspond to experimental observations exceeding the upper bounds estimated from the theoretical models. All experimental observations for E and G are lower than the upper bounds, whereas a few experimental values for ν and K exceed the upper bounds. Variability, and uncertainties in experimental determination and bound computation may explain this observation. Only some experimental values of K exceed the bounds, their numbers in the training database are identified and depicted in Figure 10(e) so that the reader can easily extract them from the data collection in case of interest.

Positions of the experimental observations with respect to the lower bounds are not shown since lower bounds being null and the elastic constants non-negative, all experimental observations within the training dataset satisfy the theoretical lower bounds.

We compare the CPU time associated with the creation of the predictor functions based on the training dataset, and the realization of one prediction (using the already created predictor functions) in Figure 11. ANN takes much longer to build the predictor functions than the other methods. Once the predictor function is created, the prediction realization is obtained in a fraction of seconds for all the methods.

.1 A.3 Leave-One-Out Cross-Validation (LOOCV)

The Leave-One-Out Cross-Validation (LOOCV) is technique to exploit the domains associated with larger prediction error or exhibits a marked non-linear behaviour ([Fuhg et al. 2020](#)). This approach consist in using a k -fold cross-validation with $k = n$ (n being total the number of observations). For each observation $i \in [1, n]$, a surrogate model \mathcal{M}_{-i} is trained on $n - 1$ observations, which

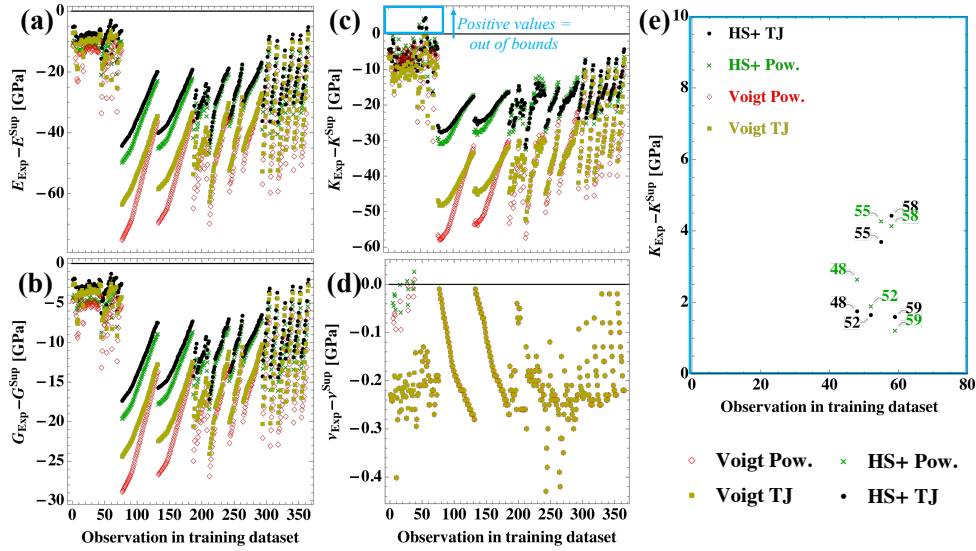


Figure 10: Difference between the experimental values in the dataset (subscript Exp) and the Voigt and Hashin-Shtrikman upper bounds (superscript sup) of (a) E , (b) G , (c) K and (d) ν . (e) Identification numbers of the few experimental observations of K exceeding the theoretical upper bounds.

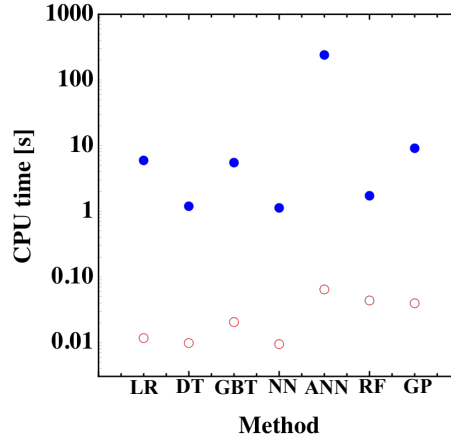


Figure 11: CPU time associated with the creation of the predictor functions based on the training dataset (blue full dots), and with one prediction using the already created predictor functions (red empty dots) for the various ML methods.

constitute a subset \mathcal{M}_{-i} (this training stage can be computationally expensive). The accuracy is finally computed using (Fuhg et al. 2020):

$$e_{LOOCV}(\mathbf{x}_i) = |\mathcal{M}(\mathbf{x}_i) - \mathcal{M}_{-i}(\mathbf{x}_i)|; \forall i \in [i, n] \quad (11)$$

where $\mathcal{M}(\mathbf{x}_i)$ is the metamodel of interest evaluated for the input \mathbf{x}_i . A small $e_{LOOCV}(\mathbf{x}_i)$ means that suppressing the observations i will not significant affect the metamodel. In other words, the interpolations made around \mathbf{x}_i are sufficiently accurate. Conversely, a large $e_{LOOCV}(\mathbf{x}_i)$ means that the information around \mathbf{x}_i is lacking.

For ANN and GP, the five observations with larger e_{LOOCV} are 5, 6, 11, 16, 17 (all from (Helmuth et al. 1966)) and 16, 17, 69, 263, 364, respectively. In the case of GP, data regarding larger w/c values is lacking. Such information can be useful to guide future experimental campaigns and optimize experiments design.

Appendix B. Hydration assemblage model

Micromechanics approaches are based on the knowledge of phases intrinsic properties and volume fractions (which evolve with time and degree of hydration). Models used to estimate the

phases fraction during hydration process are briefly exposed in this appendix.

B.1 Powers model

The ratio w/c determines the initial porosity in cement systems and can be used to estimate the porosity as a function of the degree of hydration (DOH). In the absence of filler blended in the binder, the Powers model (Powers 1960; Pichler and Hellmich 2011) allows to estimate the volume fractions of the clinker, water (capillary porosity), hydrates and chemical shrinkage (or “air”), respectively as:

$$f_{Clinker} = \frac{1 - DOH}{1 + w/c \frac{\rho_{Clinker}}{\rho_{Water}}} = \frac{20(1 - DOH)}{20 + 63w/c} \geq 0, \quad (12)$$

$$f_{Water} = \frac{\rho_{Clinker}(w/c - 0.42DOH)}{\rho_{Water} + w/c \rho_{Clinker}} = \frac{63(w/c - 0.42DOH)}{20 + 63w/c} \geq 0 \quad (13)$$

$$f_{Hydrates} = \frac{1.42\rho_{Clinker}DOH}{\rho_{Hydrates} + w/c \rho_{Clinker}/\rho_{Water}} = \frac{43.15DOH}{20 + 63w/c} \quad (14)$$

$$f_{Air} = 1 - f_{Clinker} - f_{Water} - f_{Hydrates} = \frac{3.31DOH}{20 + 63w/c} \quad (15)$$

with the mass volume of clinker $\rho_{Clinker} = 3.15 \text{ g/cm}^3$, water $\rho_{Water} = 1 \text{ g/cm}^3$ and hydrates $\rho_{Hydrates} = 2.073 \text{ g/cm}^3$ (Pichler and Hellmich 2011).

Following Hansen (Hansen 1986), the maximum degree of hydration α_{max} is a function of the w/c ratio and depends on curing conditions. For curing without external water supply, the maximum DOH denoted α_{max}^{NW} is reached when water or cement is depleted, thus

$$\alpha_{max}^{NW} = \begin{cases} \frac{w/c}{\kappa_w/\rho_{Clinker}} & \text{if } w/c \leq \kappa_w/\rho_{Clinker}, \\ 1 & \text{otherwise.} \end{cases} \quad (16)$$

For curing condition with supplementary external water supply, the maximum DOH denoted α_{max}^W is reached when cement is depleted or when the entire space available for hydrate growth, i.e. full capillary porosity, is depleted:

$$\alpha_{max}^W = \begin{cases} \frac{w/c \rho_{Clinker}}{\kappa_h - 1} & \text{if } w/c \leq (\kappa_h - 1)/\rho_{Clinker}, \\ 1 & \text{otherwise.} \end{cases} \quad (17)$$

B.2 Königsberger-Hellmich-Pichler model (KHP)

We adopt the model of the evolution of phase volume fractions proposed by Königsberger et al. (Königsberger, Hellmich, et al. 2016) propose an extension of Powers model in which C-S-H densification is accounted for, in agreement with the NMR evidence (Muller et al. 2012). C-S-H Densification is described using three hydration regimes: *regime I* dense C-S-H particle precipitated on cement particle boundaries; *regime II* C-S-H precipitates in a loosely packed configuration where gel porosity appears; and, *regime III* C-S-H precipitation completely fills the capillary porosity. The volume fractions of cement f_{cem} (approximated as clinker), other hydrates f_{CH} (assuming that portlandite CH is the main crystalline hydrate constituting the other hydration products), solid C-S-H f_{sCSH} (considered as a microporous phase with interlayer pores), gel pores f_{GP} , capillary pore f_{CP} , and void volume f_{void} (or chemical shrinkage) are, respectively, given by:

$$f_{cem} = \frac{1 - \xi}{1 + 3.185w/c} \geq 0 \quad (18)$$

$$f_{CH} = \frac{0.484\xi}{1 + 3.185w/c} \quad (19)$$

$$f_{sCSH} = \frac{1.105\xi}{1 + 3.185w/c} \quad (20)$$

$$f_{GP} = \begin{cases} 0 & 0 \leq \xi \leq \xi_{I-II} \\ \frac{4.824w/c\xi - 0.799(w/c)^2 - 0.793\xi^2}{(1+3.185w/c)(0.864w/c+1.278\xi)} & \xi_{I-II} < \xi < \xi_{II-III} \\ \frac{3.185w/c - 0.755\xi}{1+3.185w/c} & \xi_{II-III} \leq \xi \leq 1 \end{cases}$$

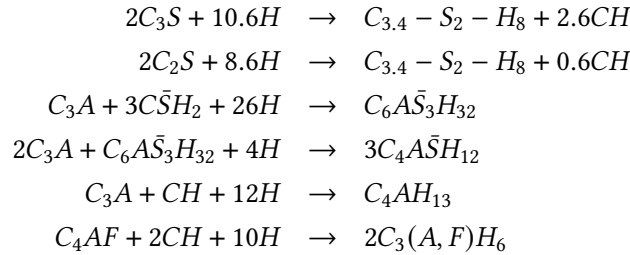
$$f_{CP} = \frac{3.185w/c - 0.755\xi}{1 + 3.185w/c} - f_{GP} \geq 0 \quad (21)$$

$$f_{void} = \frac{0.167\xi}{1 + 3.185w/c} \quad (22)$$

where $\xi_{I-II} = 0.170w/c$ and $\xi_{II-III} = 2.022w/c$ are the transition hydration degrees between hydration regimes.

B.3 Tennis and Jennings model

The phase assemblage in Tennis and Jennings (Tennis et al. 2000) model is based on the following equations:



With these stoichiometric relations and the molar volumes of the phases, it is possible to compute the volume fractions of the phases as a function of the DOH. In this approach, the aluminum bearing phases are ettringite ($C_6A\bar{S}_3H_{32}$ or Aft), monosulfoaluminate ($3C_4A\bar{S}H_{12}$ or AFm), hydrogarnet ($2C_3(A, F)H_6$) and C_4AH_{13} . With the progress of hydration, ettringite is assumed to be completely converted into monosulfoaluminate if water and C_3A are available. No phases bearing carbonates are taken into account.

This model enables to distinguish between LD and HD C-S-H, as well as gel pores. The volumes of C-S-H HD and LD are given, respectively, by:

$$V_{HD} = \frac{M_t - (M_r M_t)}{\rho_{HD}}; V_{LD} = \frac{M_r M_t}{\rho_{LD}} \quad (23)$$

where $\rho_{HD} = 1750 \text{ kg/m}^3$ and $\rho_{LD} = 1440 \text{ kg/m}^3$ are the "dried" densities of C-S-H HD and LD, respectively, as reported in (Tennis et al. 2000). The LD mass ratio with respect to the total mass of C-S-H denoted M_t and computed from the stoichiometric equations presented above, is denoted $M_r = 3.017(w/c)DOH - 1.347DOH + 0.538$. The volume of gel pore reads

$$V_{Gel \text{ Pores}} = V_{LD} - \frac{M_r M_t}{\rho_{HD}}. \quad (24)$$

Tennis and Jennings model is used to get the volume fraction of phases as a function of the degree of hydration for three different commercial cements studied in previous works. For a

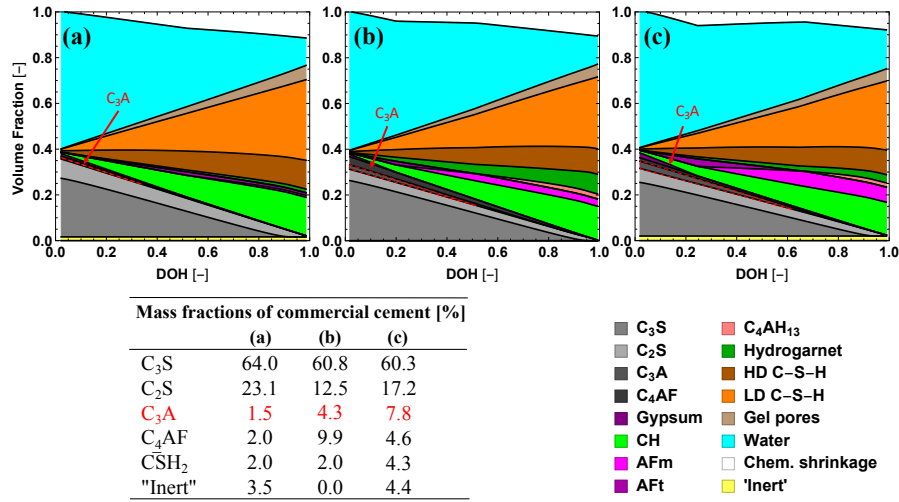


Figure 12: Volume fraction of phases for three commercial cements studied in refs. (a) (Honorio, Bary, and Benboudjema 2016; Honorio, Bary, and Benboudjema 2018; Honorio, Bary, Benboudjema, and Poyet 2016), (b) (Wyrzykowski et al. 2017) and (c) (Termkhajornkit et al. 2012), with $w/c = 0.5$, as a function of the degree of hydration (DOH). The variations on C_3A content lead to significant differences in the emergence of various Al-bearing phases.

w/c ratio equal to 0.5, the resulting phase assemblages are shown in Fig. 12. The variations on C_3A content lead to significant differences in the emergence of various Al-bearing phases: the fractions of AF-phases and C_4AH_{13} are clearly more significant in systems (b) and (c). The high- C_4AF content of cement (b) leads to a higher fraction of hydrogarnet formed.

Figure 13 shows the evolution of the volume fraction of phases with the degree of hydration as estimated using the TJ model for the samples included in the test dataset (Section 4.2.2). Note that in Figure 13(c), hydration is stopped at a DOH of approximately 0.9 since there is no water anymore available.

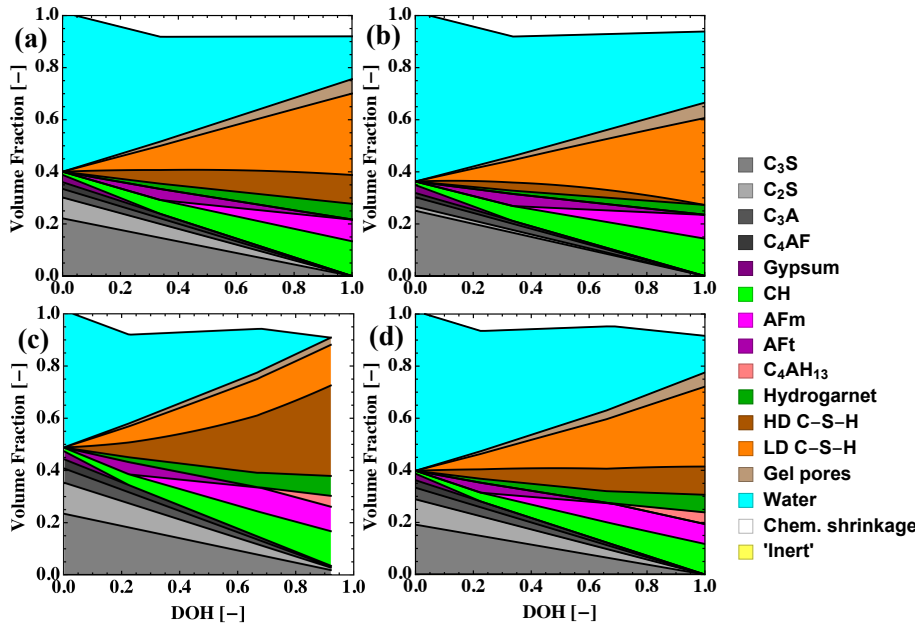


Figure 13: Evolution of the volume fraction of different phases with the degree of hydration for three commercial cements studied by (a) Constantinides and Ulm (Constantinides et al. 2004) for $w/c = 0.5$, (b) Haecker et al. (Haecker et al. 2005) for $w/c = 0.6$, (c and d) Tamtsia et al. (Tamtsia et al. 2004) for (c) $w/c = 0.35$ and (d) $w/c = 0.5$.