

New homework assignment posted to canvas

Due 2/8

Plan for today

- Mixture models
↳ overfitting
- Regression theory
- finite elements
- local polynomial reproduction
- curse of dimensionality
- intro to DNNs ↳ no black boxes

Some final comments on building probability distributions

Multivariate

Gaussian

$$\mathbf{y} \in \mathbb{R}^n$$

$$\mathbf{y} \sim N(\boldsymbol{\mu}, \Sigma) \quad \text{def}$$

$$p(\mathbf{y}) = 2\pi^{-n/2} |\Sigma| \exp \left[-\frac{1}{2} (\mathbf{y}-\boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{y}-\boldsymbol{\mu}) \right]$$

ex If $\Sigma = \text{diag}(\sigma_i^2)$ show that \mathbf{y} may be viewed as joint distribution of indep univariate gaussians

Mixture of experts

Let $C \sim \text{categorical}(\pi)$
 $\text{cat}(\pi)$

$$0 \leq \pi \leq 1$$

$$\sum_i \pi_i = 1$$

Define MOE

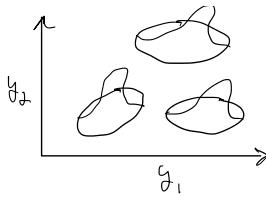
$$p(\mathbf{y} | \mathbf{x}) = \sum_c \underbrace{p(\mathbf{y} | c, \mathbf{x})}_{\text{expert model}} \underbrace{p(c | \mathbf{x})}_{\text{gating function}}$$

IDEA: Tool for Blending Multiple Models
Potentially Allowing Specialization

Example Gaussian Mixture model

$$\begin{aligned} p(\mathbf{y}) &= \sum_{c=1}^{N_c} p(c) p(\mathbf{y} | c) \\ &= \sum_{c=1}^{N_c} \pi_c N(\mathbf{y}; \boldsymbol{\mu}_c, \Sigma_c) \end{aligned}$$





Some exercises

- Compute $p(c|y)$

$$p(c|y) = \frac{p(y|c) p(c)}{p(y)} \quad \text{Bayes Rule}$$

$$= \frac{p(y|c) p(c)}{\sum_{c'} p(y|c') p(c')} \quad \text{Total Probability}$$

$$= \frac{\pi_c N(y; \mu_c, \Sigma_c)}{\sum_c \pi_{c'} N(y; \mu_{c'}, \Sigma_{c'})}$$

Remark \rightarrow Mixture models allow Bayesian inference of cluster ownership

Exercises Compute $E_{c \sim p(c|y)} [c]$
 $\text{Var}_{c \sim p(c|y)} [c]$

Overfitting in MLE

Consider fitting two clusters to scattered data w/ max. likelihood



$$\begin{aligned} \text{NLL}(\pi, \mu, \Sigma) &= \log p(y_1, \dots, y_d) = \sum_{i=1}^d \log p(y_i) \\ &= \sum_{i=1}^d \log \sum_c \pi_c p(y_i | c) \end{aligned}$$

Remark: First time we've gotten \log of sum
 \rightarrow not nice for optimization

Take $\pi_{c,i} = 1$

$$M_{C_i} = y_{d_i}$$

$$\sum_{C_i} = \varepsilon^2 I \ll 1$$

Then $NLL = -\sum_{i=1}^d \log N(y_{d_i}; y_{d_i}, \varepsilon^2 I)$

 $= -\log N(y_{d_1}; y_{d_1}, \varepsilon^2 I) - \sum_{i \neq 1}^d \log N(y_{d_i}; y_{d_i}, \varepsilon^2 I)$
 $= C + \frac{d}{2} \log \varepsilon + \frac{1}{2} \cancel{\frac{1}{\varepsilon^2}} + \text{STUFF}$

Note that $\lim_{\varepsilon \rightarrow 0} NLL = -\infty$

So that we can make the NLL arbitrarily small by attaching a cluster to a single point

Overshooting is a danger of MLE

Soln. Stabilization for GMM

$$N(y | M_C, \Sigma_C + \varepsilon^2 I)$$

$\underbrace{\hspace{10em}}$ fixed $\varepsilon > 0$
Background Noise

Another MLE example Piecewise Polynomial Regression

$$p(y|x) = \sum_c p(y|c, x) p(c|x)$$

$$p(c|x) = \text{cat}(\sigma(wx+\beta))$$

$$p(y|c,x) = N(y; A_c x + B_c, \sigma_c^2)$$

Classical Regression ~ Introductory Functional Analysis

Let $L_2(\Omega) = \left\{ f \mid \int_{\Omega} f^2 dx < \infty \right\}$

$$C_0(\Omega) = \left\{ f \mid \text{Finite Elements} \right.$$

Let $\Omega = x_0 < \dots < x_N = 1$

Let \mathcal{S} be a linear space of functions
 $\mathcal{S} = \text{span}(\{v_1, \dots, v_N\})$ satisfying

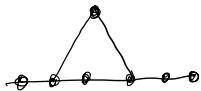
1) $v \in C^0([0, 1])$ continuous on unit interval

2) $v|_{[x_i, x_{i+1}]}^{} \text{ is a linear polynomial}$

3) $v(0) = 0$

Define $\phi_i(x)$ by Kronecker delta property

$$\phi_i(x_j) = \delta_{ij} = \begin{cases} 1 & i=j \\ 0 & \text{else} \end{cases}$$



ML def
of PL
functions

Recall $\text{ReLU}(x) = \begin{cases} x & x \geq 0 \\ 0 & \text{else} \end{cases}$

Ex Express $\phi_i(x)$ in terms of ReLU

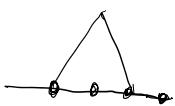
$$\phi_i(x) = A \text{ReLU}(x - x_{i-1}) + B \text{ReLU}(x - x_i) + C \text{ReLU}(x - x_{i+1})$$

$$\phi_i(x_j) = \delta_{ij} \text{ implies}$$

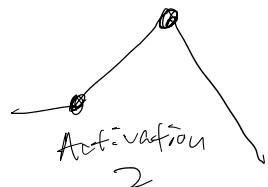
$$1 = \phi_i(x_i) = A(x_i - x_{i-1}) + 0 + 0$$

$$0 = \phi_i(x_{i+1}) = A(x_{i+1} - x_{i-1}) + B(x_{i+1} - x_i) + 0$$

$$0 = \phi_i(x_{i+2}) = A(x_{i+2} - x_{i-1}) + B(x_{i+2} - x_i) + C(x_{i+2} - x_{i+1})$$



$$A = \frac{1}{x_i - x_{i-1}} \quad B = \frac{-A(x_{i+1} - x_{i-1})}{x_{i+1} - x_i} \quad C = \frac{-A(x_{i+2} - x_{i-1}) - B(x_{i+2} - x_i)}{x_{i+2} - x_{i+1}}$$



Lemma ϕ_i form a basis for \mathcal{S}

Pf - ϕ_i is a basis if $\sum c_i \phi_i(x_j) = 0$ implies $c_i = 0$ for all i

- From Kronecker's property

$$\sum_i c_i \phi_i(x_j) = \sum_i c_i \delta_{ij}$$

Thm $\mathcal{S} \subseteq L^2([0,1])$ if $c_j = 0$

for

Define The interpolant $v_I \in \mathcal{S}$ satisfying

$$v_I(x_j) = v(x_j) \quad \text{for all } j$$

$$\text{is } v_I(x) = \sum_i v(x_i) \phi_i(x)$$

$$\begin{aligned} \text{Pf } v_I(x_j) &= \sum_i c_i \phi_i(x_j) \\ &= \sum_i c_i \delta_{ij} \\ &= c_j \end{aligned}$$

Thm Let $h = \max_i |x_{i+1} - x_i|$. If u' defined

$$\text{Then } |u(x) - u_I(x)| \leq \frac{1}{2} h^2 \max_{x \in [0,1]} |u''(x)|$$

for $x \in [x_i, x_{i+1}]$

$$\text{Pf } u(x) - u_I(x) = \int_{x_i}^x (u - u_I)'(t) dt$$

Fundamen
Thm
Calc

• Mean value thm $\Rightarrow x^* \in [x_i, x_{i+1}]$ s.t $\underset{\text{can find}}{u'(x^*)} = \frac{u(x_{i+1}) - u(x_i)}{x_{i+1} - x_i} = u'_I$

$$u'(x^*) = \frac{u(x_{i+1}) - u(x_i)}{x_{i+1} - x_i} = u'_I$$

$$\text{so } (u - u_I)'(x^*) = 0$$

ℓ

fall
of
lungs

(x^*)

$$\text{F.T.C.} \quad (u - u_I)'(t) = \int_{x^*}^t (u - u_I)''(s) ds$$

$$= \int_{x_i}^x \int_{x^*}^t (u - u_I)''(s) ds dt$$

$$= \int_{x_i}^x \int_{x^*}^t u''(s) ds dt$$

$$|u - u_I|(x) \leq \max_{s \in [x_i, x_{i+1}]} u''(s) \int_{x_i}^x \int_{x^*}^t ds dt$$

$$|u - u_I| \leq \frac{1}{2} h^2 \left(\max_{[0,1]} u'' \right)$$

✓

We can understand this informally
 \rightarrow the interpolant can match the constant + linear term of the Taylor series on each element. Integrating the remainder on the element gives h^2 scaling

Exercise Develop higher-order polynomial spaces using ReLU networks (HW problem)

First glimpse at DNN

ds

*an
alve
breath*

*c
= O*

Def A shallow neural network of width N

$$NN: \mathbb{R}^{N_{\text{in}}} \rightarrow \mathbb{R}^{N_{\text{out}}}$$

with ReLU activation consists of
a "hidden layer"

$$L_H(x) = \text{ReLU}(w x + b)$$

$$w \in \mathbb{R}^{N \times N_{\text{in}}}, \quad b \in \mathbb{R}^N$$

composed w/ a "linear layer"

$$L_L(x) = C x$$

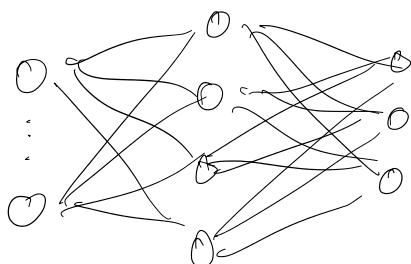
$$C \in \mathbb{R}^{N_{\text{out}} \times N}$$

$$NN_{\theta}(x) = L_L \circ L_H(x)$$

and we denote trainable params

$$\theta = \{w, b, C\}$$

We draw this
as a network w/
weights + biases
represented as
edges



crude model of brain

Some history

This is an example of a ...

Some history

This is an example of a "Multilayer Perceptron" (MLP)

More generally a deep NN / dense / feed forward can consist of many layers and other nonlinear activation functions

$$y(x) = L_1 \circ L_{H_m} \circ \dots \circ L_k(x)$$

$$L_{H_i}(x) = \sigma(w_i x + b_i)$$

activation weight bias

1958 - single hidden, trainable linear only

1970 - first backprop

2003 - renewed interest due to architectures

$$NN(x) = \sum_{i=1}^n c_i \bar{\Phi}_i(x; \theta)$$

$$\bar{\Phi}_i(x; \theta) = \sigma(w_i x + b_i)$$

Then Consider $N_{in} = N_{out} = 1$ (D scalar data)

$$\mathcal{D} = \left(x_i, y_i \right)_{i=1}^{N_D} \quad N_D = N - 1$$

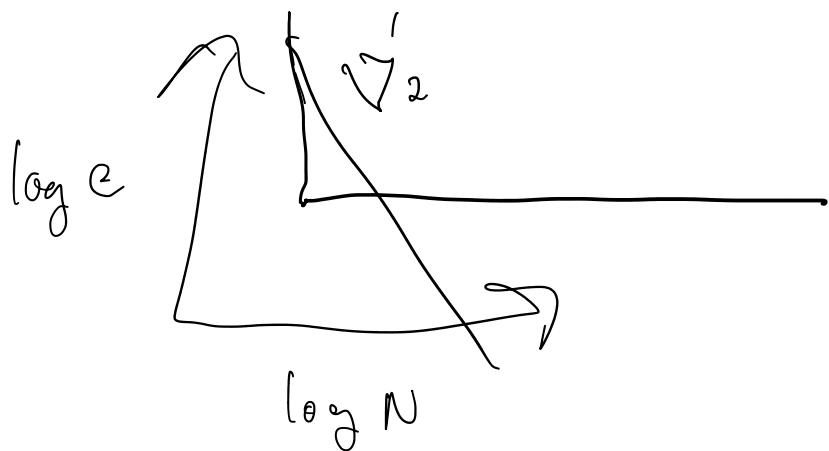
satisfying $x_i \neq x_j \forall i, j$

satisfying $x_i \neq x_j \wedge i \neq j$

then

$$\min_{\theta} \sqrt{\frac{1}{N_0} \sum_{i=1}^{N_0} |NN_{\theta}(x_i) - y_i|^2} \leq C h^2$$

PF Set: bins to match nodes
 weights to of
 linear layer to interpolant of
 data



- This is drastically oversimplified

- Does an interpolant exist for arbitrary scattered data?
- How does that depend on dimension

- A NN can evaluate PW lines in 1D — what about higher dimensions?
- When can a NN beat a traditional finite element space?
- How do we initialize a DNN when training w/ gradient descent?

Understanding the curse of dimensionality

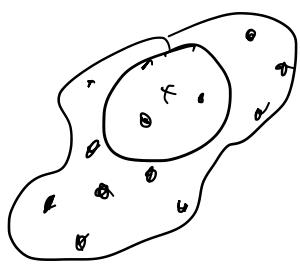
(See "Scattered data approximation"
By Holger Wendland)

- Consider a scattered dataset
$$\mathcal{S} = \{x_i, y_i\}_{i=1}^n, \quad x \in \mathbb{R}^d, y \in \mathbb{R}$$
- Aim to build a "quasi-interpolant" of \mathcal{S}
$$s(x) = \sum_{j=1}^n y_j u_j(x) \quad \text{TB}$$

where $s \approx y$ at x

- Def fill-distance

$$h_{x, \Omega} = \sup_{x \in \Omega} \min_{1 \leq j \leq n} \|x - x_j\|_2$$



the largest distance one needs to go to find a data pt-

- def A local polynomial reproduction of degree ℓ assigns u_j s.t.

polynomial reproducing

$$\textcircled{1} \quad \sum_j p_j u_j(x) = p(x) \quad \forall p \in \Pi_\ell(\mathbb{R}^d)$$

not too big

$$\textcircled{2} \quad \sum_j |u_j(x)| < c, \quad \forall x \in \Omega$$

compact support

$$\textcircled{3} \quad u_j(x) = 0 \quad \text{if } \|x - x_j\|_2 > c_2 h_{x, \Omega}, \quad \forall x \in \Omega$$

For a polynomial reproduction of order m on a compact Ω

$$|f(x) - s(x)| \leq C h_{x, \Omega}^{m+1} \|f\|_{C^{m+1}}$$

proof

$$\text{where } \|f\|_{C^{m+1}} = \max_{|\alpha|=m+1} \|D^\alpha f\|_{C^\infty}$$

Pf Let $p \in \Pi_m$ be a polynomial of order m . Then

$$|f(x) - s(x)| \leq |f(x) - p(x)| + |p(x) - \sum_j f_j u_j(x)|$$

$$\begin{aligned}
&= \|f - p\| + \left| \sum_j p_j u_j(x) - f_j u_j(x) \right| \\
&\leq |f(x) - p(x)| + \sum_j |p_j - f_j| |u_j(x)| \\
&\leq \|f - p\|_{L^\infty(B(x, c_2 h_{x, \Omega}))} \left(1 + \sum_j |u_j(x)| \right) \\
&\leq (1 + C_1) \|f - p\|_{L^\infty(B(x, c_2 h_{x, \Omega}))}
\end{aligned}$$

Taking p as the Taylor polynomial

$$\|f - p\|_{L^\infty(B(x, c_2 h_{x, \Omega}))} \leq C h_{x, \Omega}^{m+1} \|f\|_{C^{m+1}}$$

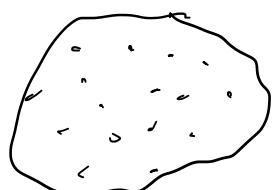
- ② So error depends on scaling of $h_{x, \Omega}$
- How does that depend on dimension?

def separation distance $g_x = \frac{1}{2} \min_{i \neq j} \|x_i - x_j\|$

def data sites X are quasi-uniform if

$$g_x < h_{x, \Omega} < C_{qu} g_x, \text{ for } C_{qu} > 0$$

informally → $g_x \sim h_{x, \Omega}$



then Given N quasi-uniform data sites

$$CN^{-1/d} \leq h_{x, \Omega} \leq CN^{-1/d}$$

$$CN^{-1/d} \leq h_{x,\Omega} \leq CN^{-1/d}$$

PF $\Rightarrow \Omega \subseteq \bigcup_{i=1}^N B(x_i, h_{x,\Omega})$

$$\text{vol}(\Omega) \leq \sum \text{vol}(B(x_i, h_{x,\Omega})) \leq CN h^d$$

$$N^{-1/d} < Ch$$

\Leftarrow Since Ω bounded, there is a x_0, R s.t

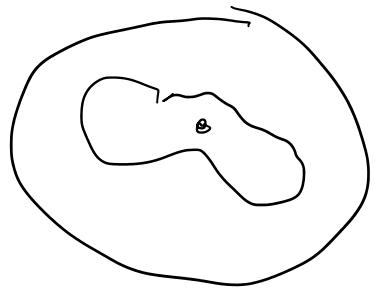
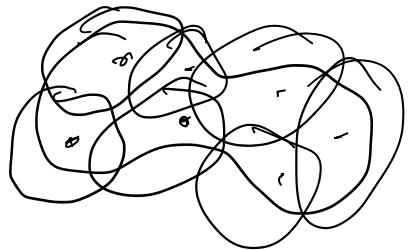
$$\Omega \subseteq B(x_0, R)$$

$\underbrace{\bigcup_{i=1}^N B(x_i, g_x)}_{\text{disjoint balls}} \subseteq \Omega$

$$N g_x^d \leq \text{vol}(\Omega)$$

$$g_x \leq CN^{1/d}$$

$\underbrace{g_x}_{\text{quasi-uniformity}} \leq CN^{-1/d}$



Punchline To hit a given error $\sim n^{-\frac{m+1}{d}}$

$$\text{Musilli} \sim C h^{m+1} \sim C N^{\frac{-m+1}{d}}$$

Today - Introduction to the finite element method and solving diff-eq. models

Control Volume Analysis

Recall the generalized Stokes theorem

Gauss Div thm

$$\int_C \nabla \cdot F dH = \int_{f \in \mathcal{E}(C)} F \cdot dA$$

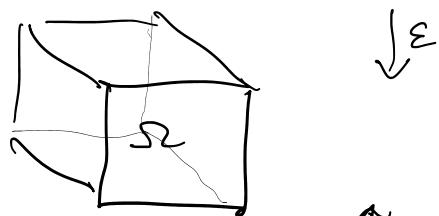
Stokes/Green's thm

$$\iint_f \nabla \times F dA = \oint_{\mathcal{E}(f)} F \cdot dl$$

Fund. Thm of calculus $\oint_e \nabla \phi \cdot dl = \phi(e_+) - \phi(e_-)$

Conservation laws

- Consider an arbitrary subdomain Ω



- Express change of system as balance of flux in/out

$$\frac{d}{dt} \int_{\Omega} \text{stuff} dx = \oint_{\partial\Omega} \text{Flux(stuff)} \cdot dA$$

ex Cons of Mass

$$\frac{d}{dt} \text{MASS} = \frac{d}{dt} \int_{\Omega} \rho dx = \int_{\partial\Omega} \rho \vec{u} \cdot dA = \frac{\text{MASS}}{\text{FLUX}}$$

Applying GDT

$$\int_{\Omega} (\partial_t \rho + \nabla \cdot \rho u) dx = 0$$

Must hold independent of choice for Ω

\Rightarrow integrand = 0

$$\partial_t \rho + \nabla \cdot \rho u = 0 \quad \text{PDE Model}$$

Examples ① Conservation of Momentum

- $\rho \partial_t u = \nabla \cdot \vec{\sigma}(u)$ (solids, u = displacement)
- $\rho \partial_t u = \nabla \cdot \vec{\sigma}(u)$ (fluids, u = velocity)

② Conservation of thermal energy

$$\rho c_p \partial_t T = \nabla \cdot \vec{q}(+) \quad q \rightarrow \text{heat flux}$$

③ Conservation of magnetic flux

$$\partial_t B = \nabla \times (v \times B) \quad v \times B \rightarrow \text{magnetic flux}$$

We will consider the abstract problem

$$\left\{ \begin{array}{l} \partial_t u + \nabla \cdot F = f(x, t) \\ F = N(u; \theta) + g(x, t) \end{array} \right. \quad \begin{array}{l} (A) \\ (B) \end{array}$$

Notes

- CLA is known physics
- CLR is unknown closure
- CLA encodes conservation structure independent of θ

$$\begin{aligned} \int_{\Omega} \partial_t u + \nabla \cdot F \, dx &= \int_{\Omega} f \, dx \\ &= \frac{d}{dt} \int_{\Omega} u \, dx + \int_{\partial\Omega} F \cdot \hat{n} \, dA = \int_{\Omega} f \, dx \\ \Rightarrow \frac{d}{dt} \int_{\Omega} u \, dx &= 0 \quad \text{if } F \cdot \hat{n} \Big|_{\partial\Omega} = 0 \\ &\quad f = 0 \end{aligned}$$

Classical Modelling

Develop $N(u; \theta)$ either
empirically or analytically

so that

① Easily calibrate θ
from experiments

e.g. rheometry

② Guarantee CL is:
well-posed

- easily solved
- computer vs analytically
- physically realizable
- amenable to UQ

Data-driven modeling

Instead use data to select F_θ
via an equality constrained opt.
problem.

$$\underset{\theta}{\operatorname{argmin}} \quad \|u - u_{\text{data}}\|^2 + \varepsilon^2 \|F - F_{\text{data}}\|^2$$

s.t. CC holds

We have the pieces:

- probability
 - regression
 - how to numerically solve PDES
- last piece \rightarrow

Next week - Finite Elements vs DNNs for PDEs

How to solve heat eqn w/ FEM

Let $\Omega = [0, 1]$

Solve $\partial_t u + c^2 \partial_{xx} u = f$

Eqn 1 $u(x=0, t) = u_L = 0$ for now

$$u(x=1, t) = u_R = 0$$

Variational Formulation

Choose trial space $u(x) \in H_0^1(\Omega)$

$$H_0^1 = \left\{ u^2 < \infty, \int u'^2 < \infty, u|_{\partial\Omega} = 0 \right\}$$

IDEA $\rightarrow u$ are candidate solutions

Pick a test space V ← more on how to choose this later

IDEA $v \in V$ will be how we measure the error

In continuous Galerkin method $V = H_0^1(\Omega)$

Variational procedure. Let $v(x) \in V$

$$\partial_t u + c^2 \partial_{xx} u = f$$

$$\int_{\Omega} v (\partial_t u + c^2 \partial_{xx} u) dx = \int_{\Omega} vf dx$$

$$\frac{d}{dt} \int_{\Omega} vu dx + c^2 \int_{\Omega} v \partial_{xx} u dx = \int_{\Omega} vf dx$$

Integrate 2nd term by parts

$$\frac{d}{dt} \int_{\Omega} vu dx - c^2 \left[\partial_x v \partial_x u + v \partial_x^2 u \right]_0^1 = \int_{\Omega} vf dx$$

$v(0) = v(1) = 0$

With the notation $(f, g) = \int fg dx$

$$\frac{d}{dt} (v, u) - c^2 (\partial_x v, \partial_x u) = (v, f)$$

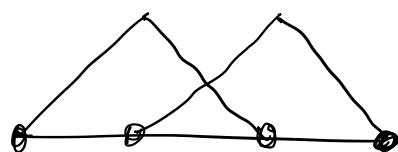
Why rewrite in this way?

- Any soln of "strong form" (Eqn 1) will satisfy this equation
- Eqn 1 needs u'' to be defined, while this one only needs to be able to integrate products of first derivatives (hence "weak form")
- exposes symmetry that we'll use to prove stuff

How to solve

Define $V_h = \{ \text{piecewise linear} \leq \text{s.t. B.C.s} \}$

$$u = \sum_{i=1}^{N-1} \hat{u}_i \phi_i(x)$$



$$= \Phi(x)^T \hat{u}$$

And substitute in

$$\frac{d}{dt} (v, u) - c^2 (\partial_x v, \partial_x u) = (v, f)$$

$$\frac{d}{dt} \sum \hat{v}_i (\phi_i, \phi_i) \hat{u}_i - c^2 \hat{v}_i (\partial_x \phi_i, \partial_x \phi_i) \hat{u}_i = \sum \hat{v}_i (\phi_i, f)$$

$$\frac{d}{dt} \sum_{i,j} \hat{v}_i (\phi_i, \phi_j) \hat{u}_j - c^2 \hat{v}_i (\partial_x \phi_i, \partial_x \phi_j) \hat{u}_j = \sum_i \hat{v}_i (\phi_i, f)$$

Regrouping

$$\sum_i \left[\hat{v}_i \left(\sum_j (\phi_i, \phi_j) \hat{u}_j - c^2 (\partial_x \phi_i, \partial_x \phi_j) \hat{u}_j - (\phi_i, f) \right) \right]$$

$\underbrace{\quad}_{M_{ij}}$ $\underbrace{\quad}_{S_{ij}}$ $b_i = 0$
mass matrix stiffness matrix

$$\hat{v}^\top \left(\frac{d}{dt} M + c^2 S \right) \hat{u} = \hat{v}^\top b$$

Require to hold for any $v \in V_h$

$$\Rightarrow \boxed{M \frac{d}{dt} \hat{u} + c^2 S \hat{u} = b}$$

linear system of ODES

To solve, discretize in time

ex if $\hat{u}^n = \hat{u}(t^n)$ explicit Euler

$$\frac{d}{dt} \hat{u}(t^n) \approx \frac{\hat{u}^{n+1} - \hat{u}^n}{\Delta t}$$

$$M \ddot{u}^{n+1} = (M - c^2 \Delta t S) \ddot{u}^n + \Delta t b^n$$

Required calculation : Build M, S, b

Gauss quadrature $\int_a^b f(x) dx = \sum_q w_q f(x_q)$

2 pt $x_1 = \frac{-1}{\sqrt{3}}, x_2 = \frac{1}{\sqrt{3}}$

or w/ change
of var

$$w_1 = w_2 = 1$$

$$\int_a^b f(x) dx = \frac{b-a}{2} \sum_i w_i f\left(\frac{b-a}{2}\xi_i + \frac{a+b}{2}\right)$$

