

New homework assignment posted to canvas

Due 2/8

Plan for
today

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

Some final comments on building probability distributions

MultivariateGaussian

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

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

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

ex If $\Sigma = \text{diag}(\sigma_i^2)$ show that 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

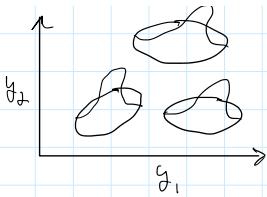
$$p(y|x) = \sum_c \underbrace{p(y|c,x)}_{\text{expert model}} \underbrace{p(c|x)}_{\text{gating function}}$$

IDEA: Tool for Blending Multiple Models
 Potentially Allowing Specialization

Example Gaussian Mixture model

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





Some exercises

- Compute $p(c|y)$

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

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

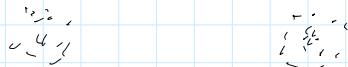
$$= \frac{\pi_c N(y; \mu_c, \Sigma_c)}{\sum_c \pi_{c'} N(y; \mu_{c'}, \Sigma_{c'})}$$

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 $\pi_c = 1$

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

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

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

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

$$= C + \frac{d}{2} \log \varepsilon + \frac{1}{2} \frac{\rho}{\varepsilon^2} + STUFF$$

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

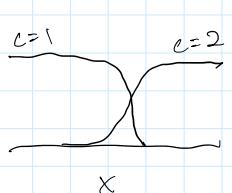
Soln. Stabilization for GMM

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

$\underbrace{\quad}_{\text{fixed } \varepsilon > 0}$
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

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

$$C_0(\Omega) = \left\{ f \mid \dots \right\}$$

Finite Elements 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

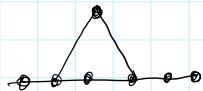
1) $v \in C^0([0, 1])$ continuous on unit interval

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

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

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

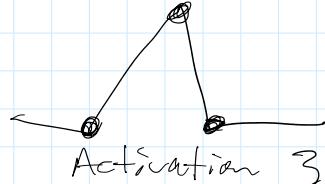
$$1 = \phi_i(x_i) = A(x_i - x_{i-1}) + 0 + 0$$

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

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



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



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

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

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

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

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

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

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

l

fa
of
lus

(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

∂S

non
alive
boreum

$I = \emptyset$

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"

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

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

composed w/ a "linear layer"

$$L_L(x) = Cx$$

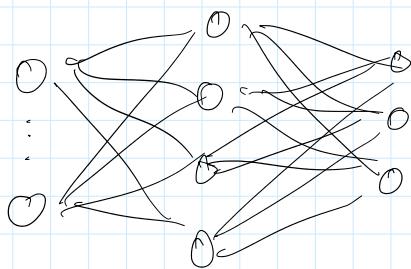
$$C \in \mathbb{R}^{N_{\text{out}} \times N}$$

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

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

$$y(x) = L_C \circ L_{H_M} \circ \dots \circ L_K(x)$$

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

activation weight bias

1958 - single hidden, trainable linear only

1970 - first backprop

2003 - renewed interest due to architectures

Then Consider $N_{in} = N_{out} = 1$

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

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

then

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

$$\leq C h^2$$

1D scalar data

PF Set: bones 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 emulate PW linear 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",)

(See "Scattered data approximation"
By Holger Wendland)

- Consider a scattered dataset

$$\mathcal{S} = \{x_i, y_i\}_{i=1}^n, \quad x_i \in \mathbb{R}^d, y_i \in \mathbb{R}$$

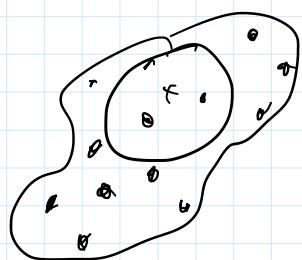
- Aim to build a "quasi-interpolant" of \mathcal{S}

$$s(x) = \sum_{j=1}^n y_j u_j(x) \quad \text{TBD}$$

where $s \approx y$ at x

- Def fill-distance

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



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

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

"not too big"

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

compact support

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

compact support

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

For a polynomial reproduction of order m
on a compact Ω

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

skip proof

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

Pf Let $p \in \mathbb{T}_m$ be a polynomial
of order m . Then

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

Taking p as the Taylor polynomial

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

④ So error depends on scaling of $h_{x,\Sigma}$

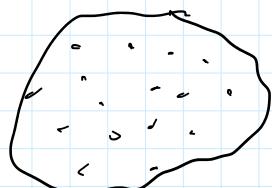
→ How does that depend on dimension?

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

def data sites X are quasi-uniform if

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

informally $\rightarrow g_x \sim h_{x,\Sigma}$



thm Given N quasi-uniform data sites

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

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

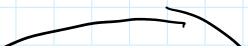
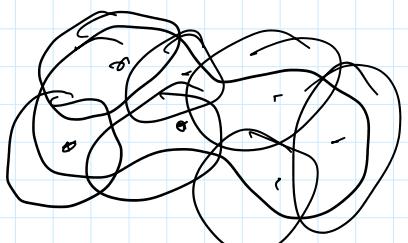
$$\text{vol}(\Sigma) \leq \sum \text{vol}(B(x_i, h_{x,\Sigma}))$$

$$\ll CN h^d$$

$$N^{-1/d} \leq Ch$$



Since Σ bounded, there
is a x_Σ, R s.t.



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

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

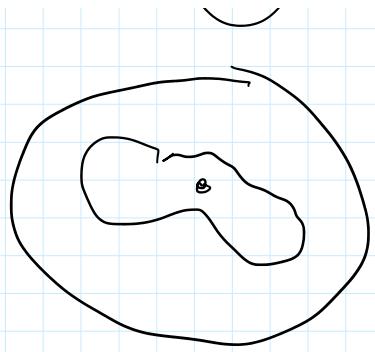
$\xrightarrow{\text{disjoint balls}}$

$$\bigcup_{i=1}^N B(x_i, g_x) \subseteq \Omega$$

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

$$g_x \leq C N^{1/d}$$

$\xrightarrow{\text{quasi-uniformity}}$ $h_{x, \Omega} < C N^{-1/d}$



Punchline To hit a given error

$$\|u - s\| \sim c h^{m+1} \sim C N^{-\frac{m+1}{d}}$$

