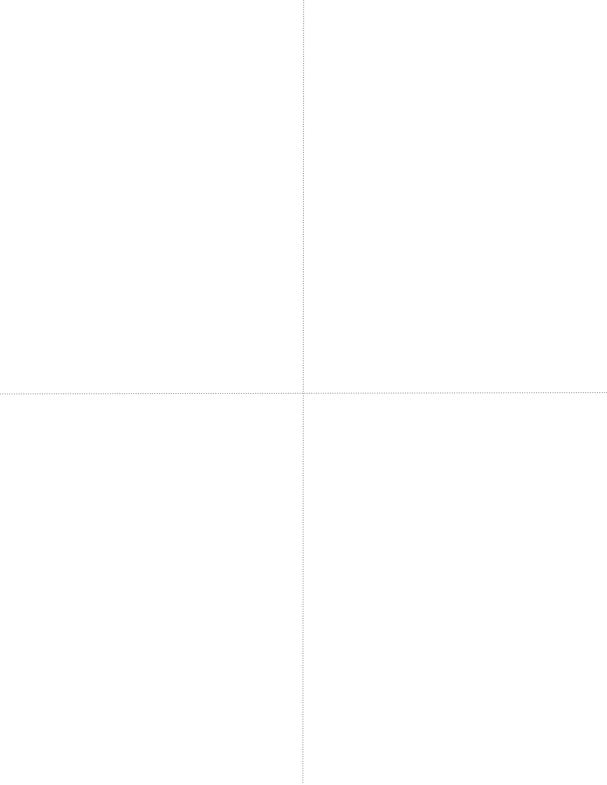
Heidelberg School

NN=FT & ARCH IS ACC
YOU NEED



ALL HITECTURE 1 " Buch Fections all you need"

OUTLINE

+T -> NN (martable!)

where ℓ is a loss function such as $\ell_{MSE} = (\phi_{\theta}(x_{\alpha}) - y_{\alpha})^2$. One may optimize θ by gradient descent

$$\frac{d\theta_i}{dt} = -\nabla_{\theta_i} \mathcal{L}[\phi_{\theta}], \qquad (11)$$

or other algorithms, e.g., classics like stochastic gradient descent (SGD) 23, 24 or Adam 25, or a more recent technique such as Energy Conserving Descent 26, 27. Throughout, t is training time of the learning algorithm unless otherwise noted.

Second, statistics. When a NN is initialized on your computer, the parameters θ are initialized as draws

$$\theta \sim P(\theta)$$
 (12)

from a distribution $P(\theta)$, where means "drawn from" in this context. Different draws of θ will give different functions ϕ_{θ} , and a priori we have no reason to prefer one over another. The prediction $\phi_{\theta}(x)$ therefore can't be fundamental! Instead, what is fundamental is the

average prediction and second moment or variance:

$$\mathbb{E}[\phi_{\theta}(x)] = \int d\theta P(\theta) \,\phi_{\theta}(x) \tag{13}$$

$$\mathbb{E}[\phi_{\theta}(x)\phi_{\theta}(y)] = \int d\theta P(\theta) \phi_{\theta}(x)\phi_{\theta}(y), \qquad (14)$$

as well as the higher moments. Expectations are across different initializations. Since we're physicists, we henceforth replace $\mathbb{E}[\cdot] = \langle \cdot \rangle$ and we remember this is a statistical expectation value. It's useful to put this in our language:

$$G^{(1)}(x) = \langle \phi_{\theta}(x) \rangle$$
 (15)

$$G^{(2)}(x, y) = \langle \phi_{\theta}(x)\phi_{\theta}(y) \rangle,$$
 (16)

the mean prediction and second moment are just the one-point and two-point correlation functions of the statistical ensemble of neural networks. Apparently ML has something to do with field theory.

Putting the dynamics and statistics together, we have an ensemble of initial θ -values, each of which is the starting point of a trajectory

The Setup.

11

Understanding ML at the very least means understanding neural networks. A neural network is a function

$$\phi_{\theta} : \mathbb{R}^d \to \mathbb{R}$$
 (1)

with parameters θ . We've chosen outputs in \mathbb{R} because, channeling

with parameters
$$\theta$$
. We've chosen outputs in \mathbb{R} because, channeling Coleman, scalars already exhibit the essentials. We'll use the lingo

Input:
$$x \in \mathbb{R}^d$$
 (2)
Output: $\phi_{\theta}(x) \in \mathbb{R}$ (3)

Network:
$$\phi_{\theta} \in \text{Maps}(\mathbb{R}^d, \mathbb{R})$$
 (4)
Data: \mathcal{D} , (5)

where the data D depends on the problem, but involves at least a subset of \mathbb{R}^d , potentially paired with labels $y \in \mathbb{R}$.

With this minimal background, let's ask our central question:

Question: What does a NN predict?

For any fixed value of θ , the answer is clear: $\phi_{\theta}(x)$. However, the answer is complicated by issues of both dynamics and statistics. First, dynamics. In ML, parameters are updated to solve problems

and we really have trajectories in Parameter Space: $\theta(t) \in \mathbb{R}^{|\theta|}$ (6)

Function Space: $\phi_{\theta(t)} \in \text{Maps}(\mathbb{R}^d, \mathbb{R})$. (8) governed by some learning dynamics determined by the optimization algorithm and the nature of the learning problem. For instance, in supervised learning we have data

Output Space: $\phi_{\theta(t)}(x) \in \mathbb{R}$

$$\mathcal{D} = \{(x_{\alpha}, y_{\alpha}) \in \mathbb{R}^d \times \mathbb{R}\}_{\alpha=1}^{|\mathcal{D}|},$$

and a loss function

)3

$$\mathcal{L}[\phi_{\theta}] = \sum_{\alpha}^{|\mathcal{D}|} \ell(\phi_{\theta}(x_{\alpha}), y_{\alpha}), \qquad (10)$$

to think of
$$\theta(t)$$
 drawn as

(7)

14

(17)

 $\theta(t) \sim P(\theta(t)),$ a density on parameters that depends on the training time and yields time-dependent correlators

 $\theta(t)$, and therefore we have an ensemble of trajectories. We choose

$$G_t^{(1)}(x) = \langle \phi_{\theta}(x) \rangle_t$$
 (18)

$$G_t^{(2)}(x,y) = \langle \phi_{\theta}(x)\phi_{\theta}(y)\rangle_t, \tag{19}$$

where the subscript t indicates time-dependence and the expectation is with respect to $P(\theta(t))$. Of course, assuming that learning is helping, we wish to take $t \to \infty$ and are interested in

$$G_{\infty}^{(1)}(x) = \text{mean prediction of } \infty\text{-number of NNs as } t \to \infty.$$

Remarkably, we will see that in a certain supervised setting there is an exact analytic solution for this quantity.

ALELACIZON ->

Question: What is a field theory? At the very least, a field theory needs

• Fields, functions from an appropriate function space, or sections

- of an appropriate bundle, more generally.
 - Correlation Functions of fields, here expressed as scalars $G^{(n)}(x_1,\ldots,x_n) = \langle \phi(x_1)\ldots\phi(x_n)\rangle.$ (128)

$$G^{(n)}(x_1, \dots, x_n) = (\phi(x_1) \dots \phi(x_n)).$$
 (128)

You might already be wanting to add more beyond these minimal requirements - we'll discuss that in a second. For now, we have

Answer: a FT is an ensemble of functions with a way to compute their correlators.

DEFINING DATA OF NN-FT (Op. P(0))

and identically distributed (i.i.d.).

]5

17

(33)

where the set of network parameters is $\theta = \{w_{ii}^{(0)}, w_i^{(1)}\}\$ independently

 $w_{ii}^{(0)} \sim P(w^{(0)})$ $w_{i}^{(1)} \sim P(w^{(1)}).$

Under this assumption, we see Observation: The network is a sum of N i.i.d. functions.

This is a function version of the Central Limit Theorem, generalizing the review in Appendix A, and gives us the Neural Network / Gaussian Process (NNGP) correspondence,

NNGP Correspondence: in the $N \to \infty$ limit, ϕ is drawn from a Gaussian Process (GP).

(34)

 $\lim_{N \to \infty} \phi(x) \sim \mathcal{N} (\mu(x), K(x, y)),$

$$\lim_{N \to \infty} \phi(x) \sim \mathcal{N}(\mu(x), K(x, y)), \qquad (34)$$

with mean and covariance (or kernel) $\mu(x)$ and K(x, y).

By the CLT, $\exp(-S[\phi])$ is Gaussian and therefore $S[\phi]$ is quadratic in networks. Now this really feels like physics, since the infinite neural network is drawn from a Gaussian density on functions, which defines a generalized free field theory. We will address generality of the NNGP correspondence momen-

tarily, but let's first get a feel for how to do computations. To facilitate then, we take $P(w^{(1)})$ to have zero mean and finite variance,

 $\langle w^{(1)}w^{(1)}\rangle = \mu_2,$ $\langle w^{(1)} \rangle = 0$ which causes the one-point function to vanish $G^{(1)}(x) = 0$. Following

Williams 34, we compute the two-point function in parameter space

(with Einstein summation) $G^{(2)}(x, y) = \frac{1}{N} \langle w_i^{(1)} \sigma(w_{ij}^{(0)} x_j) w_k^{(1)} \sigma(w_{kl}^{(0)} y_l) \rangle$ (36) $=\frac{1}{N}\langle w_i^{(1)}w_k^{(1)}\rangle\langle \sigma(w_{ij}^{(0)}x_j)\sigma(w_{kl}^{(0)}y_l)\rangle$

$$= \frac{1}{N} \langle w_i^{(1)} w_k^{(1)} \rangle \langle \sigma(w_{ij}^{(0)} x_j) \sigma(w_{kl}^{(0)} y_l) \rangle$$

= $\frac{\mu_2}{N} \langle \sigma(w_{ij}^{(0)} x_j) \sigma(w_{il}^{(0)} y_l) \rangle$,

(37)(38)

N= CHANNELS CNN

Our minimal requirements get us a partition function (129)

(131)

18

 $Z[J] = \langle e^{\int d^d x J(x)\phi(x)} \rangle$

that we can use to compute correlators, where at this stage we are

agnostic about the definition of $\langle \cdot \rangle$. In normal field theory, the $\langle \cdot \rangle$ is

defined by the Feynman path integral

 $Z[J] = \int \mathcal{D}\phi e^{-S[\phi] + \int d^dx J(x)\phi(x)},$ (130)

Euclidean Answer: a FT is a statistical ensemble of functions.

which requires specifying an action $S[\phi]$ that determines a density on functions $\exp(-S[\phi])$. But that's not the data we specify when we specify a NN. The NN data $(\phi_{\theta}, P(\theta))$ instead defines

 $Z[J] = \int d\theta P(\theta) e^{\int d^dx J(x)\phi_{\theta}(x)}$ ~ for quartum, i for question! Q: WHEN NN-QFT?

ONE A: GW SATISFY OS-AXIONS
(see 2021 paper & paper of Feeks)

FREE THEORIES (MGP consynders)

For simplicity, we again consider a single-layer fully connected network of width N, with the so-called biases turned off for simplicity: $\phi(x) = \frac{1}{\sqrt{N}} \sum_{i=1}^{N} \sum_{j=1}^{d} w_{ij}^{(1)} \sigma(w_{ij}^{(0)} x_{j}),$

 $G^{(2)}(x, y) = \mu_2 \langle \sigma(w_{ij}^{(0)} x_j) \sigma(w_{il}^{(0)} y_l) \rangle,$ where we emphasize there is now no summation on i. This is an exact-in-N two-point function that now requires only on the computation of the quantity in bra-kets. One may try to evaluate it exactly by doing the integral over $w^{(0)}$. If it can't be done, Monte Carlo estimates may be obtained from M samples of $w^{(0)} \sim P(w^{(0)})$

 $G^{(2)}(x, y) \simeq \frac{\mu_2}{M} \sum_{i=1}^{M} \sigma(w_{ij}^{(0)} x_j) \sigma(w_{il}^{(0)} y_l).$

where the last equality follows from the ones being i.i.d., $\langle w_i^{(1)} w_k^{(1)} \rangle =$

 $\mu_2 \delta_{ik}$. The sum over i gives us N copies of the same function, leaving

In typical NN settings, parameter densities are easy to sample for convenience, allowing for easy computation of the estimate. If the density is more complicated, one may always resort to Markov chains, e.g. as in lattice field theory.

With this computation in hand, we have the defining data of this NNGP. (41)

 $\lim \phi(x) \sim \mathcal{N}(0, G^{(2)}(x, y))$.

The associated action is

 $S[\phi] = \int d^dx d^dy \, \phi(x) \, G^{(2)}(x, y)^{-1} \, \phi(y),$ (42)

where $\int d^d y G^{(2)}(x, y)^{-1}G^{(2)}(y, z) = \delta^{(d)}(x - z).$ (43)defines the inverse two-point function. In fact, this allows us to

determine the action of any NNGP with $\mu(x) = G^{(1)}(x) = 0$, by computing the $G^{(2)}$ in parameter space and inverting it. So certain large neural networks are function draws from generalized free field theories. But at this point you might be asking yourself

Question: How general is the NNGP correspondence? 11= WIDTH

N= HEADS ... (le /cmg)

INTERACTIONS (in 'h) iii) INTERACTIONS CLT -> FREE INTERACTIONS = NOT FREE -> NOT CLT $\phi(x) = \frac{1}{\sqrt{N}} \sum_{i=1}^{N} \omega_i \varphi_i(x)$ of treating of UT assump. CONTROLLED INTERACTIONS $G^{(4)} = \langle \phi(x)\phi(y)\phi(z)\phi(w) \rangle$ (45) $= \sum \langle w_i w_j w_k w_l \rangle \langle \varphi_i(x) \varphi_j(y) \varphi_k(z) \varphi_l(w) \rangle$ (46)I) N FINITE IN CORRECTIONS $= \sum \langle w_i^4 \rangle \langle \varphi_i(x)\varphi_i(y)\varphi_i(z)\varphi_i(w) \rangle$ (47)+ $\sum \langle w_i^2 \rangle \langle w_j^2 \rangle \langle \varphi_i(x) \varphi_i(y) \varphi_j(z) \varphi_j(w) + \text{perms} \rangle$. G(2r) ~ 1/r-1 (48)One can see that you have to be careful with indices. The connected 4-pt function is 38 (->0 + r)2, => Gamesian thy) $G_c^{(4)}(x, y, z, w) = G^{(4)}(x, y, z, w) - (G^{(2)}(x, y)G^{(2)}(z, w) + perms),$ 2) BREAK IMPPENDENCE and watching indices carefully we obtain $G_c^{(4)}(x,y,z,w) = \frac{1}{N} \Big(\mu_4 \left\langle \varphi_i(x) \varphi_i(y) \varphi_i(z) \varphi_i(w) \right\rangle$ eg e willi -) (wiwi Ywyluj) (50) $-\mu_2^2 \left(\left\langle \varphi_i(x)\varphi_i(y) \right\rangle \left\langle \varphi_i(z)\varphi_i(w) \right\rangle + \text{perms} \right) \right),$ (51)mixing tems forbid IV) SYMMETRIES ACT, REDEFFARE WEFGHT $\left| \frac{\Phi}{\Phi} \right| \quad \mathcal{I}_{\Phi^{4}} [1] = \left\langle e^{\sqrt{1} \phi} e^{-\sqrt{1} \phi' \phi} \right\rangle_{fS}$ NW-INTERD: $\phi = \phi_{w_1,w_2,y_3,0}$ とつ => e - 50 4(x) DEFORMS P(w'w"6") $\phi(x) = \alpha \sigma(6x), \quad x \rightarrow -x, P(6) = P(-6)$ (ORE >> EPLA)PIGS [\$(-X1) -- \$(-x2)] BREAKS INDEP! EMBEDDING FORNACISM (Dirac) $\frac{\text{ASSORG}}{\text{ASS}} = \mathbb{E}_{P(a)P(-a)} \left[\phi(x,) - \phi(-x,) \right]$ SO(0+1, 1) = CFTD = LORD+1,1 NN-CFT HONOG LI NN ON IR DT', I P(6)=17-6) = (Epro) P(6) [(4,) -- (km)] V P.S. OF P.M.C. FREE SCALAR AS N->00 CFT ON RD $\phi(x) = \sqrt{\frac{2\text{vol}(B_{\Lambda}^{A})}{(2\pi)^{d}\sigma_{w_{1}}^{2}}} \frac{1}{\sqrt{\mathbf{w}^{(0)_{i}^{2}} + m^{2}}} w_{i}^{(1)} \cos \left(w_{ij}^{(0)} x_{j} + b_{i}^{(0)}\right), \quad (72)$ (QM:)) ANY COM THY HAS NO REP $w^{(1)} \sim \mathcal{N}\left(0, \frac{\sigma_{w_1}^2}{M}\right)$ $w^{(0)} \sim \text{Unif}(B_{\Lambda}^d)$ $b^{(0)} \sim \text{Unif}[-\pi, \pi],$ ii) NN(MARKOU PROLES]) -) REF. POSITIVE -> UNITARY QM iii) CLASSIC RESULTS (COMMUTATURES, HUP, FERMS & SUSY IN PROGRESK

architecture in all you need Outher

(1) UAT & FFNN

(2) Kalmagovov-arnold Nets (KANS)

(3) Aparity: MLP -> KAN

(4) plot, if the

(4) y necessary KNS & Symb reg)

(UAT & FENN

The Universal Approximation Theorem (UAT) is the first result in this direction. It states that a neural network with a single hidden layer can approximate any continuous function on a compact domain to arbitrary accuracy. More precisely, the origin version of the theorem due to Cybenko is

Theorem 2.1 (Cybenko). Let $f: \mathbb{R}^d \to \mathbb{R}$ be a continuous function on a compact set $K \subset \mathbb{R}^d$. Then for any $\epsilon > 0$ there exists a neural network with a single hidden layer of the form

$$\phi(x) = \sum_{i=1}^{N} \sum_{j=1}^{d} w_i^{(1)} \sigma(w_{ij}^{(0)} x_j + b_i^{(0)}) + b^{(1)}, \tag{21}$$

 $\theta = \{w_{ij}^{(0)}, w_i^{(1)}, b_i^{(0)}, b^{(1)}\}, \text{ where } \sigma : \mathbb{R} \to \mathbb{R} \text{ is a non-polynomial nonlinear activation function, such that}$

$$\sup_{x \in \mathcal{X}} |f(x) - \phi(x)| < \epsilon. \tag{2}$$

The architecture in Eq. [21] is known by many names, including the perceptron, a fully-connected network, or a feedforward network. The parameter N is known as the width. It may be generalized to include a depth dimension L encoding the number of compositions of affine transformations and non-linearities. Such an architecture is known as a multi-layer perceptron (MLP) or a deep feedforward network.

The UAT is a powerful result, but it has some limitations. First, though the error ϵ does get better with N, it doesn't say how many neurons are needed. Second, it doesn't say how to train the network: though there is a point θ^* in parameter space that is a good approximation to any f, existence doesn't imply that we can find it, which is a question of learning dynamics.

To ask the obvious,

Question: Why does this UAT work?

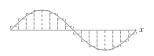


Figure 1: The Universal Approximation Theorem can be understood by approximating a function like $\sin(x)$ with a series of bumps.

In Cybenko's original work, he focused on the case that σ was the sigmoid function

$$\sigma(x) = \frac{1}{1 + e^{-x}},$$
 (23)

17

which allows us to get a picture of what's happening. The sigmoids appear in the network function of the form

$$\sigma \left(w_i^{(0)} \cdot x + b_i^{(0)}\right)$$
 (24)

which approximates a shifted step function as $w^{(0)} \to \infty$. A linear combination can turn it into a bump with approximately compact support that gets scaled by $w^{(1)}$. These bumps can then be put together to approximate any function; see, e.g., Fig. [2.1]

Cybenko's version of the UAT was just the first, and there are many generalizations, including to deeper networks, to other activation functions, and to other domains. See [28] for the original paper and [29] for a generalization to the case of multiple hidden layers.

16

Z KAN 17	NOTE:
THM: (Kolnogorev-Arnold) GIVEN MULTIVARIATE CONTIMOUS $f: [0,1]^n \longrightarrow R$	(a) Assumption of funger of Holden they come
$\exists A \text{REPRESENTATION}$ $f(x_1,, x_n) = \underbrace{\overset{\text{an}}{\Sigma}} \underbrace{\underbrace{f_2(\overset{\text{e}}{\Sigma} \varphi_{P_1}g(x_1P))}}_{g_{-0}} \underbrace{f_2(\overset{\text{e}}{\Sigma} \varphi_{P_2}g(x_1P))}_{g_{-0}} f$	(3) LOCKS LIKE NN! ACTIVATIONS ON EDGES [DEFN] (KAN) A DEPTH L f: R"0 -> R"L KOLMOGORON-APMULD NETWORK IS $f(x) = \sum_{i=1}^{n} \phi_{i-1}^{(i-1)} (\sum \phi_{i-1-2}^{(i-2)} ((\sum \phi_{i-1-2}^{(i)} (\times_{i-2}))))$ WHERE GR ACTIVS ON EDGES. NB & HAS PARAMS, LEAR NED "don't worm, take inplication generally e rules is train!"
1 MANY APPLICATIONS (2) A IMPLEMENTATION (ME Alides) (3) B-SPLINES transfle (4) = W, b(x) + w, z c, B, (x) - LOCALITY HELPS (AT. PORGET. - SCALING THEOREM (2.1) N = # PARAMS L= TEST RMSE L~ N d= DATA DIM poly spline THEN \(\times B-SPLINE-KAN = k+1) COMPARE: \(\times Relu-M = \frac{k+1}{d} \)	SPARSITY KAN = SPARSE MLP (20) SPARSITY MAD: () HAS GUARANTEE () MAPS ACTIVS PROM MODES > EDGES "MLP BIG" WIOTH BN2 51,, 53 (w) where 6 deep (20) (w) where 6 deep (Mat()) (w) we hat (RN2, RN2) (w) AB & def (Mat()) (v) AB & def (Mat())

() () () () () () () () () ()	3 SYMBOLS
	METHOD GIVEN DATA SYMBOLIC REGRESSION TO FIND INTERPADE FORS FOR IT.
"MLP SMALL" (middle Slatin = KAN indices) RESHAPE b -> lgk RN2 NRN = \(\gamma \begin{array}{c} \left(\ell + i) \) \[\sigma \begin{array}{c} \left(\sigma \left \right) \] = \(\sigma \begin{array}{c} \left(\ell + i) \) \[\sigma \left(\left(\ell + i) \) \] = \(\sigma \begin{array}{c} \left(\ell + i) \) \[\sigma \left(\left(\ell + i) \) \] = \(\sigma \begin{array}{c} \left(\ell + i) \) \[\sigma \left(\left(\ell + i) \) \] = \(\sigma \begin{array}{c} \left(\ell + i) \) \[\sigma \left(\left(\ell + i) \) \[\sigma \left(\left(\ell + i) \) \] = \(\sigma \left(\left(\ell + i) \) \[\sigma \lef	Qs: WHAT FORM? 1D? HIGHER D? LLM-LEX VISUAL SYMB REG
$= \sum_{\substack{\ell \neq k}} \mathcal{B}_{i\elljh}^{(\ell+i)} \sigma_{i} \left(\sum_{j=i}^{N} A_{jjh} \right) z_{j,k}^{(\ell)} $ $= \sum_{\substack{\ell \neq k}} \mathcal{B}_{iijh}^{(\ell+i)} \sigma_{i} \left(A_{ijkh} z_{k,n}^{(\hat{I})} \right) =: \sum_{k} \phi_{iq}^{(\ell+i)} (z_{n,i}^{(\hat{I})})$	CAN LM GUESS FUNCTIONAL FORM? YES, BUT ITER W GA
(varying of for MLP become "bains" for learned activation Q: MLP SMALL -> KAN SPARSETY	Prompt LLM Ansatz Fit Params Score
US RAW)OM UNIFORM SPARSITY?	Figure 2: The structure of LLM-LEX. (result: competitive of Mathematica, as of open source LMs
LIMITATION: 1) FUNCTIONS (b/c traved or)	
NOTE: ORIG KAN PACKAGE HAS ZUAGES!	
Prune Suggested? Fit KAN Prune Suggested? No Fit edges with LLM-LEx Simplify Expression Figure 6: The general structure of KAN-LEx.	

MCP SMALL

$$N_{g} = BN^{2}$$

$$RN^{2}$$

$$RN$$

$$MLP \leq MACC \qquad 2 \stackrel{\text{lt}}{\stackrel{?}{\sim}} = \stackrel{BN^2}{\stackrel{>}{\sim}} B_{(b)}^{(l+1)} = \stackrel{\text{b mod B}}{\stackrel{>}{\sim}} \left(\stackrel{\sim}{\stackrel{>}{\sim}} A_{b}^{(l)} \stackrel{\sim}{>} \stackrel{\sim}{>} \right)$$