

TAM 598 Lecture 24 :

Neural Networks

Announcements:

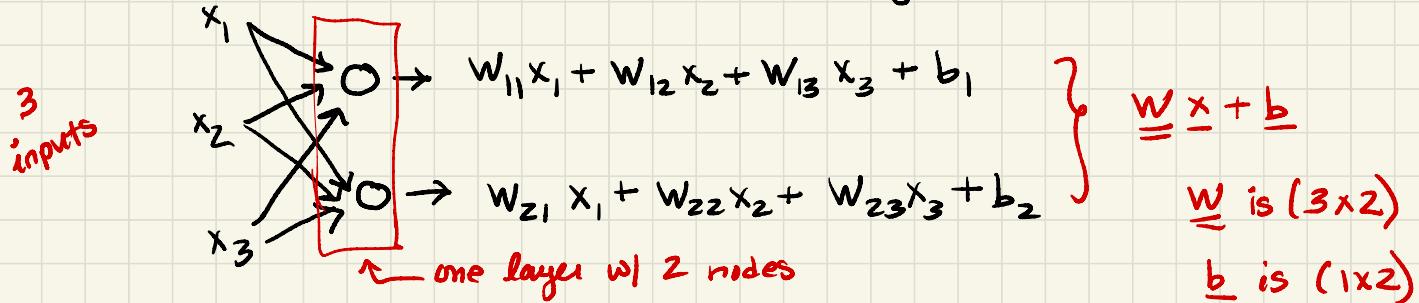
- HW b covers lectures 21-23; due on Fri May 2

Neural Networks - massively parametrized function approximators
that express hierarchically layered information (each successive
layer captures increasingly abstract features from the input)

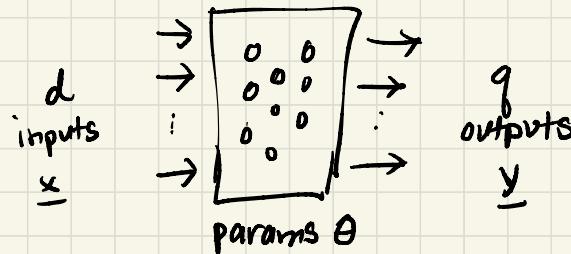
→ operate in layers

→ each layer takes as input the output of the
previous layer, and

- ① applies a linear transformation
- ② adds a constant bias
- ③ passes result through a non-linear function
- ④ passes to the next layer



"deep" \Rightarrow many layers (≥ 5 ish)



$$\underline{\underline{y}} = f(\underline{\underline{x}}; \theta)$$

parameters

one layer: $\underline{y}_1 = h^{(1)}(\underline{\underline{W}}^{(1)} \underline{x} + \underline{b}^{(1)})$

$$\underline{y}_1 = f_1(\underline{x})$$

two layers: $\underline{y}_2 = h^{(2)}\left(\underline{\underline{W}}^{(2)} \underline{y}_1 + \underline{b}^{(2)}\right)$

$$y_2 = f_2(y_1)$$

$$= f_2(f_1(\underline{x}))$$

$$= f_2 \circ f_1(\underline{x})$$

⋮

L layers $\underline{y}_n = (f_L \circ f_{L-1} \dots \circ f_1)(\underline{x})$

"fully connected
DNN"

$h^{(i)}$ is an activation

a nonlinear function applied to each argument
(ReLU, sigmoid, hyperbolic tangent, sinusoid, step function)

$\underline{w}^{(i)}$ is a weight matrix

$\underline{b}^{(i)}$ is a bias

other archs:
recurrent
LNNs
autorec
residual NNs

Most typically all but the last layer have the same activation function. The last depends on constraints of the final output, eg

1) a real number

$$\xrightarrow{\quad} \textcircled{O} \rightarrow h^{(i)}(z) = z$$

2) pmf on K categories

$$h^{(i)}(\underline{z}) = \frac{\exp(z_i)}{\sum_{j=1}^K \exp(z_j)}$$

$\left. \begin{matrix} 0 \\ 0 \\ \vdots \\ 0 \end{matrix} \right\} K \text{ nodes}$

Universal approximation theorem - given a good activation function, and enough nodes, we can approximate ANY continuous function arbitrarily well

Training Regression Networks - Loss Function View

$$\begin{array}{l} \text{inputs } \underline{x}_{1:n} = (x_1, \dots, x_n) \\ \text{outputs } \underline{y}_{1:n} = (y_1, \dots, y_n) \end{array} \quad \left. \right\}$$

$$y = f(\underline{x}; \theta)$$

need to fit
these params
(weights/biases)

$$\text{Loss Function } L(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i; \theta))^2$$

$$\text{Find } \theta^* = \operatorname{argmin} L(\theta)$$

solution: not unique, not linear, not convex,

\curvearrowright local min is
a global min
tools of alg
guaranteed
to converge

Training Regression Networks - Probabilistic View

likelihood $p(y_{1:n} | \underline{x}_{1:n}, \theta)$

maximize log likelihood $L(\theta) = -\log p(y_{1:n} | \underline{x}_{1:n}, \theta)$

same result as before if observations are independent and measurement noise is gaussian w/ mean given by the NW and constant variance

ie say $p(y_i | \underline{x}_i, \theta) = N(y_i | f(\underline{x}_i; \theta), \sigma^2)$

$$= \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{(y_i - f(\underline{x}_i; \theta))^2}{2\sigma^2} \right\}$$

from independence $p(y_{1:n} | \underline{x}_{1:n}, \theta) = \prod_{i=1}^n p(y_i | \underline{x}_i, \theta)$

So we want to minimize

$$L(\theta) = -\log p(y_{1:n} | \underline{x}_{1:n}, \theta)$$

$$= -\sum_{i=1}^n \log p(y_i | \underline{x}_i, \theta)$$

some!

$$= \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - f(\underline{x}_i, \theta))^2 + \text{const}$$

Stochastic Optimization \Rightarrow subsample available data
(in batches)

issues: $L(\theta)$ is non convex
summation in $L(\theta)$ could be huge

minim. \mathbb{E}_z [exp. over
of $L(\theta; z)$]

$$\min_{\theta} \mathbb{E}_z [l(\theta; z)]$$

z = random vector

going to
be idx of
data pt

eg let $I \sim \text{Categorical}(\frac{1}{n}, \dots, \frac{1}{n})$

pick w/ equal
prob the index of
one if the n obs.

then $\ell(\theta, I) = (y_I - f(x_I; \theta))^2$ so $Z = I$

and $E_I [\ell(\theta, I)] = \sum_{i=1}^n P(I=i) \ell(\theta, i)$

$$= \sum_{i=1}^n (y_i - f(x_i; \theta))^2$$

$$= L(\theta)$$

$E_I [\ell(\theta, I)]$
minimizing is same as minimizing $L(\theta)$

e.g w/ batch of size m

let I_1, I_2, \dots, I_m all be cat ($\frac{1}{n}, \dots, \frac{1}{n}$)

$$\text{then } l_m(\theta; I_{1:m}) = \frac{1}{m} \sum_{j=1}^m (y_{I_j} - f(x_{I_j}; \theta))^2$$

$$\text{and } \mathbb{E}[l_m(\theta; I_m)] = \mathbb{E}\left[\frac{1}{m} \sum_{j=1}^m (y_{I_j} - f(x_{I_j}; \theta))^2\right]$$

$$= \frac{1}{m} \sum_{j=1}^m \left[\mathbb{E} (y_{I_j} - f(x_{I_j}; \theta))^2 \right]$$

$$= \frac{1}{m} \sum_{j=1}^m L(\theta)$$

$$= L(\theta)$$

Stochastic Gradient Descent (aka Robbins - Monro algorithm)

want $\min_{\theta} \mathbb{E}_Z [l(\theta; Z)]$

algorithm:

- ① initialize $\theta = \theta_0$

- ② iterate $\theta_{t+1} = \theta_t - \alpha_t \nabla_{\theta} l(\theta_t, z_t)$ where
 z_t are independent samples of Z

noisy gradient

α_t is the learning rate

as long as

$$\sum_{t=1}^{\infty} \alpha_t = +\infty$$

$$\sum_{t=1}^{\infty} \alpha_t^2 < +\infty$$

} then SGD converges to
a local min of $\mathbb{E} [l(\theta, Z)]$

e.g. $\alpha_t = \frac{A}{(Bt+C)^p}$ with $p \in (0.5, 1)$

so with batches of size m at step t:

choose m indices randomly $i_{t1}, i_{t2}, \dots, i_{tm}$

$$\theta_{t+1} = \theta_t - \alpha_t \nabla_{\theta} \frac{1}{m} \sum_{j=1}^m (y_{it_j} - f(\underline{x}_{it_j}; \theta))^2$$

$$= \theta_t - 2 \alpha_t \frac{1}{m} \sum_{j=1}^m (y_{it_j} - f(\underline{x}_{it_j}; \theta)) \nabla_{\theta} f(\underline{x}_{it_j}; \theta)$$

We get the gradient $\nabla_{\theta} f$ using back-prop / autodiff

Other stochastic opt alg's: SGD w/ momentum

AdaGrad

Adam (adaptive moment estimation)