

TAM 598 Lecture 24 :

Neural Networks

Announcements:

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

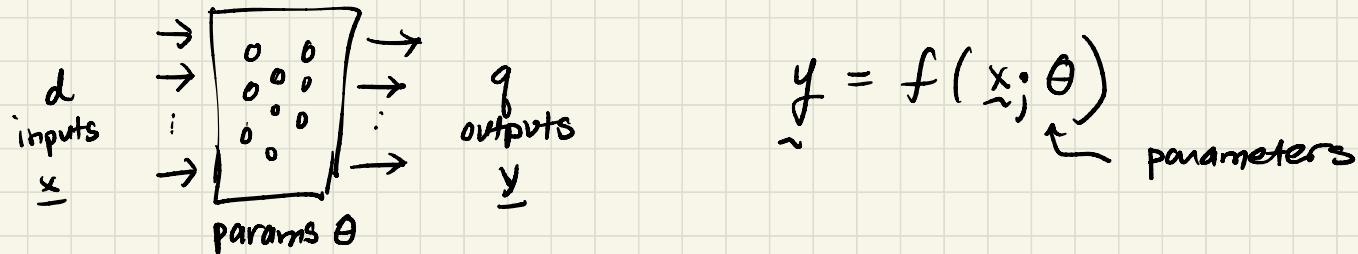
Neural Networks - massively parametrized function approximators
that express hierarchically layered information (each successive
layer captures increasingly abstract features from the input)

→ 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)



one layer :

two layers :

L layers

$h^{(i)}$ 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

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

$$\begin{array}{ccc} \searrow & & \\ & O & \rightarrow h^{(i)}(z) = z \\ \nearrow & & \end{array}$$

2) pmf on K categories

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

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

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

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

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

So we want to minimize

$$\begin{aligned} L(\theta) &= -\log p(y_{1:n} | \underline{x}_{1:n}, \theta) \\ &= -\sum_{i=1}^n \log p(y_i | \underline{x}_i, \theta) \\ &= \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - f(\underline{x}_i, \theta))^2 + \text{const} \end{aligned}$$

Stochastic Optimization

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

then $\ell(\theta, I) =$

and $E_I [\ell(\theta, I)] =$

e.g w/ batch of size m

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

then $l_m(\theta; I_{1:m}) =$

and $E[l_m(\theta; I_m)] =$

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

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)