

TAM 598 Lecture 25 :

Neural Networks, cont'd

Announcements:

- HW 7 covers lectures 24-26 ; due on Fri May 16th

Classification using Neural Networks

features $\underline{x}_{1:n}$, discrete targets $y_{1:n}$, N.N. $f(\underline{x}; \theta)$

I) Binary Classification $y = 0$ or $y = 1$

single point

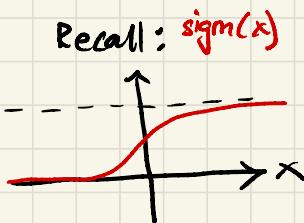
$$p(y=1 | \underline{x}, \theta) = \text{sigm}(f(\underline{x}; \theta)) = \frac{e^{f(\underline{x}; \theta)}}{1 + e^{f(\underline{x}; \theta)}}$$

$$p(y=0 | \underline{x}, \theta) = 1 - \text{sigm}(f(\underline{x}; \theta))$$

likelihood $p(y | \underline{x}, \theta) = [\text{sigm}(f(\underline{x}; \theta))]^y [1 - \text{sigm}(f(\underline{x}; \theta))]^{1-y}$

all points likelihood

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



train network by maximizing log likelihood, ie minimizing the cross entropy loss

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

$$= -\sum_{i=1}^n \left\{ y_i \log \text{sigm}(f(x_i; \theta)) + (1-y_i) \log [1 - \text{sigm}(f(x_i; \theta))] \right\}$$

II) Multiclass Classification $y = 0, 1, \dots, K-1$ K values

need: $p(y=k | \underline{x})$

use: NN $f(\underline{x}; \theta)$ with K outputs

$$f(\underline{x}; \theta) = (f_0(\underline{x}; \theta), \dots, f_{K-1}(\underline{x}; \theta))$$

K scalars

construct a probability vector:

$$P(y=k|x, \theta) = \text{softmax}_k(f(x; \theta)) = \frac{\exp(f_k(x; \theta))}{\sum_{k'=0}^{K-1} \exp(f_{k'}(x; \theta))}$$

Likelihood of full dataset

$$P(y_{1:n} | x_{1:n}, \theta) = \prod_{i=1}^n P(y_i | x, \theta)$$

Minimize the softmax cross entropy loss

$$\begin{aligned} L(\theta) &= -\log P(y_{1:n} | x_{1:n}, \theta) \\ &= -\sum_{i=1}^n \log [\text{softmax}_{y_i}(f(x_i; \theta))] \end{aligned}$$

III) Regularization

since NN's are highly flexible \Rightarrow prone to overfitting
often minimize

$$J(\theta) = L(\theta) + \lambda R(\theta)$$

positive param \downarrow \uparrow penalizes big weights, e.g.

$$R(\theta) = \sum_j \theta_j^2$$

IV) Bayesian interpretation of regularization - regularization can be seen as a maximum posterior approach, given prior $p(\theta)$

$$p(\theta | x_{1:n}, y_{1:n}) \propto p(y_{1:n} | x_{1:n}, \theta) p(\theta)$$

Maximize this

equivalent to minimize loss

$$J(\theta) = \underbrace{-\log p(y_{1:n} | x_{1:n}, \theta)}_{L(\theta)} - \underbrace{\log p(\theta)}_{\lambda R(\theta)}$$

examples:
CNNs →
process grid
like data.
Used in image
recog etc.

extract spatial
features.

apply filters
that detect
edges, shapes...

e.g.) prior $p(\theta) =$ gaussian w/ mean $\mu = 0$,
variance $1/\alpha$

$$= N(\theta | 0, \alpha^{-1} \mathbb{I})$$

then $-\log p(\theta) = \frac{\alpha}{2\pi} \sum_j \theta_j^2 + \text{const}$

i.e.) L2
regularizer