

MATHEMATICAL FOUNDATIONS OF REINFORCEMENT LEARNING

LECTURE III: APPROXIMATION ARCHITECTURES & OPTIMIZATION

ADDITIONAL REFS: "DEEP LEARNING: AN INTRODUCTION FOR APPLIED MATHEMATICIANS", BY C. HIGHAM AND D. HIGHAM, SIAM REVIEW 61(4)(2019): 860-891

TODAY:

* ARCHITECTURES

* ANN'S

* LEAST SQUARES

* NONLINEAR OPTIM.

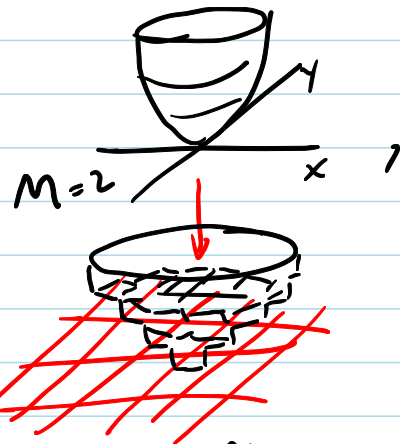
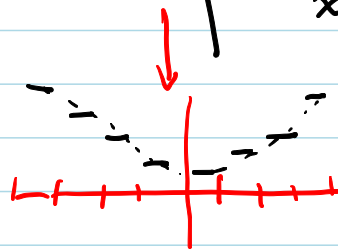
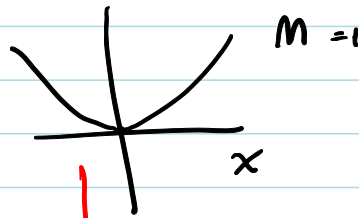
"OPTIMIZATION METHODS FOR LARGE-SCALE MACHINE LEARNING", L. BOTTOU, F.E. CURTIS

AND J. NACEDAL, SIAM REVIEW 60(2)(2018): 223-311.

BASIC APPROXIMATION ARCHITECTURES

$$J^*(x) : \mathbb{R}^m \rightarrow \mathbb{R}$$

J_F $m \leq 3$



PIEWISE CONSTANT APPROX.
 (POLYNOMIAL)

$$J^*(x) \sim \tilde{J}(x, \pi) = \sum_{l=1}^m \pi_l \phi_l(x), \text{ WHERE}$$

$$\phi_l(x) = \begin{cases} 1 & \text{if } x \in S_l, \\ 0 & \text{if } x \notin S_l. \end{cases}$$

$\underbrace{S_1, S_2, S_3}$

↳ 1 PARAMETER π_l PER SEGMENT $\rightarrow N$ SEGMENTS IN 1D

CURSE OF DIMENSIONALITY $\rightarrow N^2$ SEGMENTS IN 2D

(ONLY PRACTICAL FOR $d \leq 3$) $\rightarrow N^m$ % (N M)

AN ALTERNATIVE: GLOBAL POLYNOMIAL APPROXIMATION

SUPPOSE $x := (x_1, \dots, x_m)$, AND WE LOOK FOR A

QUADRATIC APPROXIMATION OF $J^*(x)$. DEFINE THE **BASIS**

$\phi_0(x) = 1$, $\phi_i(x) = x_i$, $\phi_{ij}(x) = x_i x_j$, THEN

$$J^*(x) \sim \tilde{J}_2(x, \pi) = \pi_0 \phi_0(x) + \sum_{i=1}^m \pi_i \phi_i(x) + \sum_{i=1}^m \sum_{j=1}^m \pi_{ij} \phi_{ij}(x)$$

+ : WE NEED $1 + m + m^2$ PARAMETERS
 $\ll N^m$

- : WE NEED TO INCREASE POLYNOMIAL DEGREE TO IMPROVE ACCURACY (TAYLOR)

$$J^*(x) \sim \tilde{J}_3(x, \pi) = \tilde{J}_2(x, \pi) + \sum_{i=1}^m \sum_{j=1}^m \sum_{k=1}^m \pi_{ijk} \phi_{ijk}(x)$$

$\hookrightarrow \phi_{ijk}(x) = x_i x_j x_k$

$1 + m + m^2 + m^3 \Rightarrow O(m^{\text{deg}})$

FOR EXAMPLE, APPROX. A 10-DIMENSIONAL J^* , $n = 10$

* PIECEWISE CONSTANT APPROX WITH 8 NODES PER DIMENSION!

8^{10} NODES $\approx 10^9$ PARAMETERS (MEMORY?)

* GLOBAL POLYNOMIAL OF DEGREE 8:

10^8 PARAMETERS \rightarrow 1 ORDER OF MAGNITUDE LESS.
(GOOD FOR $n \approx 10$)

BOTH PIECEWISE AND GLOBAL POLYNOMIAL APPROXIMATIONS ARE

LINEAR IN π . HOW DO WE FIND BEST π 'S?

SUPPOSE WE HAVE A DATASET $\{x_i, J^*(x_i)\}_{i=1}^{N_s}$ WITH N_s

SAMPLES OF THE OPTIMAL COST: INTERPOLATION OR REGRESSION

$$\pi^* = \underset{\pi}{\operatorname{argmin}} \sum_{i=1}^{N_s} \|J^*(x_i) - \tilde{J}(i, \pi)\|^2$$

$\hookrightarrow \sum \pi_i \phi_i(x_i)$

LINEAR LEAST SQUARES! (NORMAL EQNS., SVD...)

ARTIFICIAL NEURAL NETWORKS

WE NEED APPROXIMATION ARCHITECTURES TO HANDLE $n \gg \gg 1$

\Rightarrow NONLINEAR APPROXIMATION / NONLINEAR OPTIMIZATION

1-LAYER (SHALLOW) NETWORK:

$$\tilde{J}(x, v, \pi) = \sum_{l=1}^m \pi_l \phi_l(x, v) = \langle \pi, \phi(x, v) \rangle$$

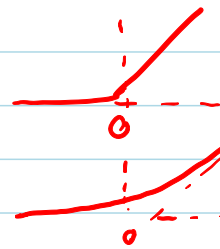
where $\phi_l(x, v) = \sigma((Ay(x) + b)_l) \rightarrow$ coordinate

ACTIVATION FUNCTION:

$$\sigma(z) = \max\{0, z\}$$

$$\sigma(z) = \ln(1 + e^z)$$

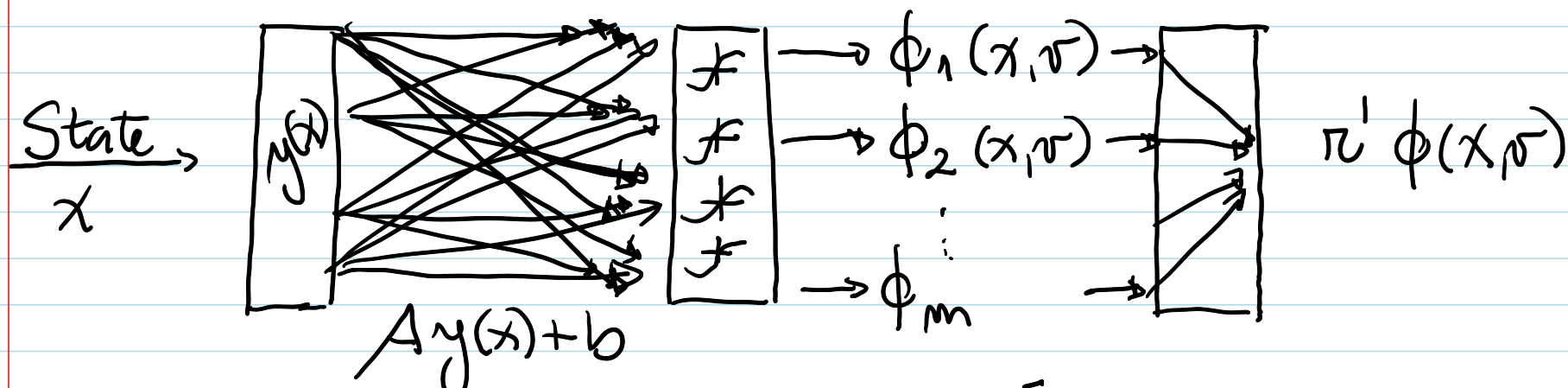
$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



ENCODER
 $y: \mathbb{R}^n \rightarrow \mathbb{R}^m$
 NONLINEAR

$\in \mathbb{R}^{m \times m}$
 $\in \mathbb{R}^m$

ENCODER - LINEAR - NONLINEAR \rightarrow LINEAR WEIGHTS.



DEFINING $\mathcal{L}^2(x) := \sigma(A^2 x + b^2)$, COMPONENTWISE

MULTILAYER (DEEP) NEURAL NETWORK

$$\tilde{J}(x, \sigma) = \mathcal{L}^{\sigma_4} \circ \dots \circ \mathcal{L}^{\sigma_3} \circ \mathcal{L}^{\sigma_2} \circ \mathcal{L}^{\sigma_1}(y(x))$$

PARAMETERS \rightarrow # OF LAYERS (DEPTH)
 \rightarrow # OF NEURONS (WIDTH)
 \rightarrow CHOICE OF σ 's, A 's, b 's
 \rightarrow ARCHITECTURE (FFN, RESNET, DNN)

TRAINING NEURAL NETWORKS

Notation! $\tilde{J}(x, \nu, \tau), \tilde{J}(x, \nu) \rightarrow \tilde{J}(x, \tau)$

GIVEN N_S SAMPLES $\{x_i, \tilde{J}^*(x_i)\}_{i=1}^{N_S}$, WE SOLVE

$$\tau^* = \underset{\tau}{\operatorname{argmin}} \sum_{i=1}^{N_S} \frac{1}{2} \|\tilde{J}^*(x_i) - \tilde{J}(x_i, \tau)\|^2$$

$:= \|g_i(\tau)\|^2$

HOWEVER, $\tilde{J}(x, \tau)$ IS NONLINEAR IN τ .

\Rightarrow NONLINEAR, NON-CONVEX, LARGE-SCALE OPTIMIZATION.

NEED TO COMPUTE OPTIMALITY CONDITIONS REQUIRING $\nabla_{\tau} g_i(\tau)$, USE CHAIN RULE/BACKPROPAGATION.

In PARTICULAR,

$$\frac{\partial \|g_r(\tau)\|^2}{\partial A^v(i_{ij})} = -g_r(\tau) \underbrace{\left[\sum_{m+1} \right]}_{A \times b} \underbrace{\left[\sum_m \right]}_{\begin{bmatrix} \sigma \\ \vdots \\ \sigma \end{bmatrix}} \underbrace{\left[\sum_{m-1} \right]}_{\begin{bmatrix} \sigma \\ \vdots \\ \sigma \end{bmatrix}} \underbrace{\left[\sum_{v+1} \right]}_{\sum_{v+1} I_{ij}} \underbrace{\left[\sum_v \right]}_{\sum_v I_{ij}} \underbrace{\left[\sum_{r-1} \right]}_{\sum_{r-1} \dots}$$

OPTIMALITY CONDITIONS

$$\min_{\tau \in \mathbb{R}^m} f(\tau)$$

τ^* IS A LOCAL MINIMUM IF $\exists \epsilon > 0$ SUCH THAT

$$f(\tau^*) \leq f(\tau), \quad \forall \tau \text{ WITH } \|\tau - \tau^*\| < \epsilon$$

τ^* IS A GLOBAL MINIMUM IF $f(\tau^*) \leq f(\tau) \quad \forall \tau \in \mathbb{R}^m$

NECESSARY OPTIMALITY CONDITIONS:

LET π^* BE A LOCAL MIN, $g \in C^1$. THEN $\nabla g(\pi^*) = 0$
(1st-ORDER), AND $\nabla^2 g(\pi^*)$ IS POSITIVE SEMIDEFINITE. (2nd)

SUFFICIENT OPTIMALITY CONDITIONS:

$\nabla g(\pi^*) = 0$ AND $\nabla^2 g(\pi^*)$ POSITIVE DEFINITE

$\Rightarrow \pi^*$ IS A STRICT LOCAL MIN OF g ,

$\exists \alpha > 0$ AND $\epsilon > 0$ SUCH THAT $g(\pi) \geq g(\pi^*) + \alpha \|\pi - \pi^*\|^2$
FOR $\|\pi - \pi^*\| < \epsilon$.

EXERCISE: GO THROUGH THE PROBLEM

$$\min_{\pi \in \mathbb{R}^m} \frac{1}{2} \pi^t Q \pi - b^t \pi, \quad Q \text{ SYMMETRIC } m \times m$$

GRADIENT METHODS:

Solving $\min_{\pi \in \mathbb{R}^m} g(\pi)$ THROUGH THE

SEQUENCE $\pi_{t+1} = \pi_t + \gamma_t s_t$, $t=0, \dots$,

WHERE $\gamma_t > 0 \forall t$, AND s_t IS A **DESCENT**

DIRECTION $\nabla g(\pi_t)^t s_t < 0$, UNTIL

$\nabla g(\pi_t) = 0$, OR SUFFICIENTLY SMALL.

How to choose (γ_t, s_t) ?

GRADIENT METHODS!

CHOOSING S_t :

i) STEEPEST DESCENT: $S_t = -\nabla g(r_t)$

ii) NEWTON'S METHOD: $S_t = -(\nabla^2 g(r_t))^{-1} \nabla g(r_t)$

(IDEA: TO MINIMIZE A QUADRATIC APPROX OF g AROUND r_t)

iii) QUASI-NEWTON METHOD: $S_t = -D_t \nabla g(r_t)$

$D_t > 0$, SYMMETRIC (EMULATES $\nabla^2 g(r_t)$)

STEPSIZE RULES: $\gamma_t = C$, OR $\gamma_t \xrightarrow{t \rightarrow \infty} 0$
↓
SMALL

CONVERGENCE:

- ASSUME $\exists c_1, c_2 > 0$ SUCH THAT $\forall t$

$$c_1 \|\nabla g(\pi_t)\|^2 \leq -\nabla g(\pi_t)' S_t$$

$$\|S_t\| \leq c_2 \|\nabla g(\pi_t)\|$$

$$\pi_{t+1} = \pi_t + \gamma_t S_t$$

- ASSUME $L > 0$ $\|\nabla g(\pi) - \nabla g(\bar{\pi})\| \leq L \|\pi - \bar{\pi}\|$

AND i) $0 < \gamma < 2c_1 / (Lc_2^2)$

ii) $\gamma_t \rightarrow 0$, $\sum_{t=0}^{+\infty} \gamma_t = \infty$

$\Rightarrow g(\pi_t) \rightarrow -\infty$ OR $\lim_{t \rightarrow \infty} \nabla g(\pi_t) = 0$.

VARIANTS:

STEEPEST DESCENT WITH MOMENTUM

$$\pi_{t+1} = \pi_t - \gamma \nabla g(\pi_t) + \beta (\pi_t - \pi_{t-1}) ; \gamma > 0, 0 \leq \beta < 1$$

OR

$$\pi_{t+1} = \pi_t - \gamma \sum_{k=0}^t \beta^{t-k} \nabla g(\pi_k)$$

// PROXIMAL POINT ALGORITHM: (IMPLICIT)

$$\pi_{t+1} = \pi_t - \gamma \nabla g(\pi_{t+1})$$

⇓

$$\pi_{t+1} = \text{PROX}_{\gamma g}(\pi_t) := \underset{\pi}{\text{argmin}} \left(g(\pi) + \frac{1}{2\gamma} \|\pi - \pi_t\|^2 \right)$$

STOCHASTIC GRADIENT DESCENT

MOTIVATION: $\min_{\pi} f(\pi) := \frac{1}{N_S} \sum_{k=1}^{N_S} f_k(\pi) \stackrel{!}{=} \frac{1}{N_S} \sum_{k=1}^{N_S} \|\nabla f_k(\pi)\|^2$

WHAT IF THE TRAINING SET N_S IS TOO LARGE?

→ NEED TO DO 1 BACKPROPAGATION / GRADIENT FOR EACH SAMPLE, AT EACH ITERATION (EXPENSIVE)

LET $K(t)$ BE SEQUENCE OF INDEPENDENT UNIFORM VARIABLES OVER $\{1, \dots, N_S\}$

SGD: $\pi_{t+1} = \pi_t - \gamma_t \nabla f_{K(t)}(\pi_t)$, WITH $\gamma_t > 0$.

OR

SGD (BATCHES): $\pi_{t+1} = \pi_t - \underbrace{\gamma_t}_{\#S(t)} \sum_{k \in S(t)} \nabla f_k(\pi_t) \stackrel{!}{=} \frac{1}{N_S} \sum_{k=1}^{N_S} \nabla f_k(\pi_t)$

SUMMARY:

- * LINEAR APPROXIMATION FOR $J^*(i)$
- * ARTIFICIAL NEURAL NETWORKS
- * TRAINING ANN'S
- * GRADIENT DESCENT METHODS
- * STOCHASTIC GRADIENT DESCENT