

PERCEPTIONS (3)

- CLASS, PAC, VC (12)
- INFER, LOG-LIKELIHOOD (4)
- CROSS-ENTR
- MLP (6)
- GRAD. DESCENT (15)
- AUTO. DIFF. (12)
- REGULARIZATION (13)
- CNN (12)
- ATTENTION, NLP (28)
- RL (8)

PERCEPTIONS

- McCulloch-Pitts 43'
- FROM LOGIC GATES TO CONNECTIONIST MODEL OF BRAIN (IMAGE RECOGNITION)
- $x \in \mathbb{R}^n$ IS INPUT
- $w \in \mathbb{R}^n$ WEIGHTS
- b THRESHOLD (OR BIAS: $b = -g$)
- $\hat{y}_w(\bar{x}) = \text{sign}(\bar{w} \cdot \bar{x} + b)$
- w, b IS THE FULL MODEL
- NOT BIOLOGICAL:
 - NO TIME VARIABLE (INSTANT RESPONSES)
 - UNCONSTRAINED INPUT (INSTEAD OF \mathbb{R}^n)
 - SUMMATION OF INPUTS IS LINEAR AS IN DOT PRODUCT

INFEERENCE

- $S \sim \text{iid}$, MODEL $D(y|x)$: $y \sim P(\theta = g)$
- $P: X \rightarrow [0, 1]^{|Y|}$: $y \sim D_x$ (NORMALISED)
- $H = \{P_{\theta}\}_{\theta}$
- PERCEPTION: $P_y(\bar{x}) = (1 + e^{-\bar{w} \cdot \bar{x} + b})^{-1}$
- AKA. LOGISTIC REGRESSION (HERE, A GREATER θ PRODUCES MORE PEAKED RESULTS)
- CLASSIFIER:
 - STOCHASTIC (SAMPLE FROM DISTRIBUTION)
 - DETERMINISTIC (TAKE θ)
- ADDING BACK A PRIOR IS EASY, ADD $-\frac{1}{m} \log P(\theta)$ (m^{-1} ENURES NORMALISATION)
- BUT CHOICE OF PRIOR IS NOT OBVIOUS
- BAYES-OPTIMAL CLASSIFIER IF THE CLASSIFIER HAS TO BE DETERMINISTIC, USE $h_g(\bar{x}) = \arg\max_{\theta} P_{\theta}(\bar{x})$ yet THIS IS BECAUSE THE MAXIMUM DISTRIBUTION IS CONCENTRATED AT THE MAXIMUM.
- ENSEMBLING: AVERAGE $P(\bar{x})$ OF MORE MODELS (ALL WITH LOW $|x_{\text{out}}$) → $h_{\text{opt}}(\bar{x}) = \arg\max \int d\theta P(\theta) e^{-\ln(S(\bar{x}))} P_{\theta}(\bar{x})$ THEORETICALLY MINIMISES $E[\text{error}]$

CLASSIFICATION

PAC-LEARNING FORM

ASSUMPTION $\Rightarrow \forall \epsilon, \delta \in (0, 1)^2$, IF CONDITION, $\Pr[\text{error} < \epsilon] \geq 1 - \delta$

RISK (DETERMINISTIC) (AKA. LOSS → $\Pr[\text{error}]$)

$$R(h) = \Pr_{\bar{x} \sim D_x} [h(x) \neq f(x)] = \mathbb{E}_{\bar{x} \sim D_x} [1_{h(x) \neq f(x)}]$$

EMPIRICAL RISK (ER) SAMPLE AVERAGE OF RISK

$$R_S(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{h(x_i) \neq y_i}$$

FREQUENTIST APPROACH: AN ER ALGORITHM PICKS $h_S \in H$ ST. ER IS MINIMISED

$$\mathcal{L}: (X \times Y)^m \rightarrow H$$

PAC-LEARNABLE $X \sim D_x(y)$, y IS ASSIGNED DETERMINISTICALLY BY NATURE

D_{X,F} REALIZABLE IN H

$\exists h^* \in H$ ST. $R_{D_x,F}(h^*) = 0$ (UP TO $\Rightarrow R_S(h^*) = 0$ a.s.)

IF REALIZABLE, M ≥ M_{H,E,S} ⇒ $\Pr_{S \sim D^m} [R_{D_x,F}(h(S)) \leq \epsilon] \geq 1 - \delta$ THE λ TAKES ONLY A M DATASETS

H FINITE IS PAC-REALIZABLE → $\lambda = \text{ERM}_H$

$$M_H(E, \delta) = \left\lceil \frac{\ln(1/\delta)}{\epsilon^2} \right\rceil$$

COROLLARS: IF REALIZABILITY KNOWS, $\forall \epsilon \in (0, 1)$, $\Pr_{S \sim D^m} [R_{D_x,F}(h(S)) \leq \epsilon] \geq 1 - \delta$ (THE BOUND CAN BE VACUOUS, ADDING NOTHING)

VC-DIMENSION OF H IS THE MAX SIZE OF C SET THAT CAN BE SHATTERED BY H:

- VC_{dim} IS LARGER THAN THE SIZE OF ANY SET THAT CAN BE SHATTERED BY H
- PROVING VC_{dim} INVOLVES FINDING THE UPPER BOUND ON SHATTERED SETS

PAC-LEARNABLE H CLASS IFF VC_{dim} IS FINITE, THEN ANY ER ALG IS PAC:

$$W_H(E, \delta) = \frac{1}{\epsilon^2} \cdot \frac{d \log(1/\delta)}{\epsilon^2}$$

PERCEPTRONS ON N-DIM. DATA HAVE N+1 VC_{dim}

RISK $R_D(h) = \Pr_{\bar{x} \sim D_x} [h(x) \neq y]$ (REDUCES TO PREVIOUS CASE IF DETERMINISTIC)

BAYES-OPTIMAL CLASSIFIER $D(x, y) = D_x(x) D(y|x)$

MINI-BATCHES ARE FED IN RANDOM ORDER AND SHUFFLED. SGD IS JUST A NOISY GRADIENT ESTIMATE SO IT CAN BE COMBINED WITH ALL OTHER METHODS

REDUCE INDUCTIVE BIAS: KNOWING M, GOAL E, S PICK H ST. VC_{dim}(H) = d < ∞ ANY ERM IS SUCCESSFUL

REDUCE INDUCTIVE BIAS: KNOWING M, GOAL E, S PICK H ST. VC_{dim}(H) = d < ∞ ANY ERM IS SUCCESSFUL

MULTIUSER PERCEPTION (MLP)

UNIVERSAL APPROXIMATION THM:

- K COMPACT
- $g: K \rightarrow \mathbb{R}$, C⁰
- $\sigma: \mathbb{R} \rightarrow \mathbb{R}$, C⁰, POLYNOMIAL OR BOUNDED
- MLP HAVE VC_{dim} < ∞ (RELU AGN-PAC-LEARN)
- UAT SAYS E_{app} CAN BE LOWERED ARBITRARILY

NETWORKS AS UNIVERSAL FITTING DEVICES

MLP ARCHITECTURE COMPOSITION OF PERCEPTIONS. L LAYERS, WHERE EACH LAYER MAPS $\mathbb{R}^{N_{l-1}} \rightarrow \mathbb{R}^{N_l}$

SOFTMAX (LAST LAYER) $P_y(\bar{x}^L) = \frac{e^{\bar{x}^L}}{\sum_{j=1}^{N_l} e^{\bar{x}^L_j}}$ NORMALISES, THEN USE CE LOSS

MATRIX NOTATION $W \bar{x}^{L-1} + b = \bar{z}^L$ (PRE-activations) $\bar{x}^L = \begin{cases} \sigma(\bar{z}^L) & l < L \\ \bar{x}^L & l = L \end{cases}$

ACTIVATION FUNCTION THERE ARE TRANSFORMATIONS OF THE PARAMETERS THAT PRESERVE THE FUNCTION

- AFFINE → WHOLE NETWORK REDUCES TO AFFINE
- SIGN (LIKE tanh) AS AN APPROXIMATION OF SIGN (AKA. STEP FUNCTION)
- RECTIFIED LINEAR UNITS (ReLU) $\text{ReLU}(x) = \max(0, x)$ (SOLVES THE VANISHING GRADIENT PROBLEM)
- GELU AS A SMOOTH APPROXIMATION

SIMMETRIES THERE ARE TRANSFORMATIONS OF THE PARAMETERS THAT PRESERVE THE FUNCTION

- DISCRETE: PERMUTING INDEXES OF UNITS. $(T_{l=1}^{L-1} (N_l)!) \text{ eq. classes}$
- CONTINUOUS: STUFF LIKE t_c AT THE LAST LAYER, OR DEPENDING ON ACTIVATION: $\alpha \otimes \text{ReLU}(ax) = a \text{ReLU}(x)$ SO INVERSE SMOOTHING CAN BE INCLUDED IN W, b

PERFORMANCE GPU HARDWARE IS OPTIMISED FOR MATRIX MULTIPLICATION, YOU CAN ALSO COMPUTE MULTIPLE MULTIPLICATIONS AT ONCE BY EXTENDING / DUPLICATING THE MATRIX. BROADCASTING $\alpha(t)$ IS ALSO EMBARRASSINGLY PARALLEL.

STOCHASTIC GRADIENT DESCENT SAMPLE THE GRADIENT FROM A SET, INSTEAD OF INDUCES RANDOMNESS: TRADEOFF BETWEEN STABILITY AND SPEED

SGD HYPERPARAMETERS CROSS-VALIDATION CAN BE USED TO FIND THE BEST HYPER-PARAMETERS:

- SPLIT DATA IN k (A "FOLD")
- TRAIN k-1 TIMES, LEAVING OUT A FOLD, WHICH IS USED FOR VALIDATION
- AVERAGE VAL. ERR. OVER k → CROSS-VAL. ERR.

EXPLORE THE RANGE OF HYPERPARAM (EX. GRID SEARCH)

K SHOULD BALANCE SPEED AND STABILITY. USING THE WHOLE TEST DATA IS RISKY

GRADIENT DESCENT

ASSUMPTIONS

- f CONTINUOUS, NBOUNDED
- LOSS EC²
- μ, σ

ITERATIVE ALGORITHM:

- SAMPLE $\bar{x}^t = \bar{x}^t + \alpha_t \nabla f(\bar{x}^t)$ TO THE VERTEX.
- SECOND ORDER → USES TAN² PARABOLOID, SET \bar{x}^t TO THE VERTEX.
- CONVERGES QUICKLY IN CONVEX SETTINGS.
- COSTS: $O(n^2)$ FOR H, $O(n^3)$ FOR MATRIX INVERSION

GRADIENT DESCENT

$\bar{x}^t = -\alpha_t \nabla f(\bar{x}^t) | \alpha_t > 0$ APPROXIMATE f NEAR \bar{x}^t AS THE TANGENT

$\alpha_t < \frac{1}{L}$ WILL DECREASE OBJECTIVE LOSS (GENERALLY DECREASES IN t)

ASSUMING f GLOBAL MINIMIZER

ISSUES:

- OVERSHOOTING FOR $\alpha > L^{-1}$ (MAY CAUSE OSCILLATIONS)
- SLOW CONVERGENCE ALONG SOME DIRECTIONS DUE TO GENTLE SLOPE

GRADIENT DESCENT

$\bar{x}^t = -\alpha_t \nabla f(\bar{x}^t) + \beta_t (\bar{x}^{t-1} - \bar{x}^t)$ LOOK-AHEAD TERM COMPENSATES FOR MOMENTUM, REMOVING OVERSHOOT.

$\beta_t \in [0.8, 0.9]$ TYPICALLY ADDS OSCILLATIONS BUT MAY OVERSHOOT. BY ADDING OSCILLATIONS AT END EASY TO IMPLEMENT

EXponential ROLLING AVERAGE ADAPT (ADAPTIVE PER-VARIABLE) RATE BASED ON LANDSCAPE.

$\alpha_t = \frac{1 - \beta}{1 - \beta^t} z^t + \beta \frac{1 - \beta^{t-1}}{1 - \beta^t} \bar{x}^{t-1}$ FAST, BUT NOISY DEPENDING ON DATA, COULD BE WORSE THAN NESTEROV.

$\bar{x}^t = \frac{1 - \beta}{1 - \beta^t} z^t + \beta \bar{x}^{t-1}$ \bar{x}^t ARE ROLLING EXP. AVERAGES OF GRADIENT, MOMENTUM

SGD HYPERPARAMETERS CROSS-VALIDATION CAN BE USED TO FIND THE BEST HYPER-PARAMETERS:

- SPLIT DATA IN k (A "FOLD")
- TRAIN k-1 TIMES, LEAVING OUT A FOLD, WHICH IS USED FOR VALIDATION
- AVERAGE VAL. ERR. OVER k → CROSS-VAL. ERR.

EXPLORE THE RANGE OF HYPERPARAM (EX. GRID SEARCH)

K SHOULD BALANCE SPEED AND STABILITY. USING THE WHOLE TEST DATA IS RISKY

INITIALISATION RELEVANT TO NON-CONVEX PROBLEMS

IN MLP SYMMETRIES WILL MAKE THE GRADIENT THE SAME $\Rightarrow \text{UNI}(0, 1)$ OR $\text{UNI}(-1, 1)$, THE VARIANCE OF THE GRADIENT IS $\text{UNI}(-\frac{1}{2}, \frac{1}{2})$

STOPPING CRITERIA SGD → IF $\|\nabla f(\bar{x}^t)\| < \epsilon$ FOR ϵ CHOSEN BEFORE

LEARNING RATE SCHEDULE YOU KNOW IT SHOULD DECREASE OVER TIME

- STEP DECAY: $\alpha \rightarrow \frac{\alpha}{m} \rightarrow \frac{\alpha}{m^2} \dots \rightarrow \frac{\alpha}{m^5}$ (TRANSITIONS FROM EXPLORATION TO EXPLOITATION)
- COSINE DECAY: $\alpha(1 + \cos(\pi \frac{t}{T}))$ CONTINUOUS DECREASES SPEND TIME ON HIGH AND SLOW DOWN TOWARDS END
- KAIMING INITIALISATION: $W_{ij} \sim \text{N}(0, \frac{2}{N_{l-1}})$ ONE OR VERY FEW PEAKS → WANTS $\alpha \approx \frac{2}{N_{l-1}}$

MINI-BATCH SIZE (SMALLER BATCH → NOISIER)

ALSO TOO LARGE MB IS TOO STABLE IN SGD MADE MORE EFFICIENT BY STACKING THE \bar{x} INTO A SINGLE, LARGER MATRIX SO WITH GPU LARGER BATCH SIZES ARE PREFERRED

BRIEF WARMUP: INCREASE α FROM 0 TO α BEFORE STARTING SOME SCHEDULER

AUTOMATIC DIFFERENTIATION

DIFFERENTIAL PROGRAMMING:
WRITE A PROGRAM WITH FREE PARAMETERS, DIFFERENTIATE THE OUTPUT WRT PARAM.
(PYTORCH → JAX)

SYMBOLIC DIFF
COMPUTER ALGORITHM THAT TAKES THE SYSTEM (SAS) AND EXPRESSIONS (SUBSTITUTION RULES)
• SOLUTION MAY BE INEFFICIENT
• CANNOT HANDLE ARBITRARY LOGIC (IF/ELSE, LOOPS)

REVERSE DIFF
COMPUTE THE ADJOINT OF EACH w_i ; VAR: $\frac{\partial w_i}{\partial w_j}$
GO BACKWARDS, UP TO GRADIENT

FORWARD DIFF
1. STATIC SINGLE ASSIGNMENT (SSA)
 $\cos(x_1^2 + x_2^3) = x_3$
 $w_1 = x_1$
 $w_2 = x_2$
 $w_3 = w_1^2$
 $w_4 = w_2^3$
 $w_5 = w_3 + w_4$
 $w_6 = \cos(w_5)$

THE DIFFERENTIATED GRAPH HAS THE SAME LOGO AS THE FUNCTION
THE ALGORITHM GROWS IN #VARS, RETAINING EFFICIENCY IN LARGE DIMENSIONAL OUTPUTS.

DUAL NUMBERS
 $\{0, \epsilon\} \cong \mathbb{R}^2$
THE TAYLOR SERIES CAN BE USED TO EXTEND FUNCTION DOMAINS TO DUAL NUMBERS
 $\frac{\partial f}{\partial x} = 0$

USING OVERLOADING IN PROGRAMMING TO COMPUTE DERIVATIVE IN ONE STEP
 $P(x) \text{ POLYNOMIAL HAS } P(x+\epsilon) = P(x) + \epsilon P'(x)$
 $f(x+\epsilon) = f(x) + \epsilon f'(x)$

FINITE DIFFERENCE FOR f BLACK BOX, APPROXIMATE 2F USING $f(x+\epsilon) - f(x)$, (MAY CAUSE FLOAT ERRORS, SINCE ϵ NEEDS TO BE $\ll 1$)

CONTROL FLOW FORWARD DIFF WORKS ON GENERAL CODE (EVEN WITH IF BLOCKS)

WORKS ON CODE ITSELF: TAKES THE FUNCTION DEFINITION AND CONSTRUCTS THE CODE OF THE DERIVATIVE, APPENDIN

BACKPROPAGATION ($\delta \in \mathbb{R}^n$, $n = \sum_{l=1}^L (N_{l-1} + 1) N_l$)
SINCE $f: \mathbb{R}^n \rightarrow \mathbb{R}^2$, IT MAKES SENSE TO USE REVERSE DIFF
(BY ASSUMPTION, $f^2 = \sigma^2(z) \circ f^1$)
 $f^1: \mathbb{R}^2 \rightarrow \mathbb{R}^L$, $L = \sum_{l=1}^L (N_{l-1} + 1) N_l$

VANISHING GRADIENTS
IF ACTIVATION IS A SATURATING FUNCTION, THE ERROR IS NOT PROPAGATED TO THE FIRST LAYERS (SLOWS LEARNING)

SOLUTION TO USE ReLU

COMPUTATION COMPUTING GRADIENTS IS ALSO GROUPED IN MATRICES FOR EFFICIENCY

DATA AUGMENTATION USE ASSUMPTIONS ON TRAINING DATA TO EXPAND TRAIN SET

ASSUME SMOOTHNESS OF d , SO IF YOU PERTURB x , $d(x+\epsilon) \approx d(x)$ (THEN JOU SAMPLE ϵ AND PRODUCE $x+\epsilon$)

EXPERIENCE (SHOW PAST) REPLAY (EXAMPLES) TO AVOID CATASTROPHIC FORGETTING (BREAKING PAST) MEMORIES

REGULARIZATION

GOAL
ONE WANTS TO MINIMIZE GENERALIZATION ERROR, TRAINING ERROR IS AN INTERMEDIATE STEP.

JUST LOWERING TRAIN E. LEADS TO OVERFITTING

THE LOSS LANDSCAPE HAS MANY LOCAL MINIMA OF SIMILAR VALUE BUT DIFFERENT GEN. ERROR.

2X PASSES ARE NEEDED:

1. FW: COMPUTE w_i VALUES

2. BW: COMPUTE THE GRADIENT

$O(m)$ SINCE YOU NEED TO COMPUTE $f(x) \times m$ TIMES

↑ YOU CAN ALSO COMBINE REVERSE/ FORWARD DIFF DEPENDING ON WHAT IS MORE TRACTABLE

(REVERSE DIFF) TAPES

"TAPES" RECORD ALL THE INTERMEDIATE VALUES AND THE COMPUTATIONAL GRAPH.

(FW → RECORD DATA)

(BW → USE DATA)

COMPUTING A OPERATION SAVES THE VALUE BUT UPDATES THE TAPE WITH VARIABLES USED TO COMPUTE THE RESULT (THIS ALSO EXTENDS TO CONTROL FLOW)

TAPES WORK AS LONG AS EVERY STEP IS DOCUMENTED IN A NEW OBJECT

(REVERSE DIFF) SOURCE-TO-SOURCE

WORKS ON CODE ITSELF:

TAKES THE FUNCTION DEFINITION AND CONSTRUCTS THE CODE OF THE DERIVATIVE, APPENDIN

LABEL Smoothing

"REAV", THE ONE-HOT LABEL VECTORS, SO THEY DON'T PRODUCE JUST PEAKED PROBABILITIES.

$q_j = (1-\epsilon) \delta_j + \frac{\epsilon}{J-1}$

BATCH NORM (CNN >)

VARIANCE OF ACTIVATION PROPAGATES EXPONENTIALLY WRT # LAYERS

(A SIMILAR THING HAPPENS IN THE BACKPROPAGATION OF ERROR)

LOGIT

ADD A TERM $+ 2/\pi^2 \cdot \log(\text{softmax}(z))$ TO LAST LAYER IS SMALLER

WON'T LEAD TO THE EXACT SOFT LABELING STRUCTURE

NORMALISE M, O^2

PER-BATCH, AT A LAYER ($z = \frac{x-\mu}{\sigma}$)

EPOCHS TAKE LONGER BUT NETWORK CONVERGES FASTER

(AVOIDS THE NEED FOR THE AW TO NORMALISE FEATURES)

LAYER NORM

IF A LAYER IS REUSED, NORMALIZATION HAPPENS BY PATTERN.

(PERSISTENCE IN PATTERN)

IN PATTERN, μ AND σ ARE SHARED ACROSS LAYERS

NORMALIZATION ACTS BASICALLY LIKE "GUIDING" INPUT

VALUES USING A FUNNEL

CORRELATED TO BETTER GENERALIZATION (OVERPARAM. NETWORKS)

NO CONSENSUS YET ON BEST WAY TO BIAS ALGORITHMS

FINAL SIZE $\sim 10^4$ BUT TOKENS WILL BE MEANINGFUL

DROPOUT

ENSEMBLING: TRAIN MORE MODELS AND PICK MAJORITY VOTE

HISTAKES SHOULD BE REFLCTD/UNCORR. COSTLY

IN FEEDFORWARD, ADD A DROPOUT LAYER $R^m \rightarrow R^n$ WHICH IN TRAINING IS $\frac{1-p}{1-p}$ WITH $p = 1-p$

DISTRIB. OF. IN VAE, $p = 1-p$, $\text{OTAKES } p$

SIMULATES E. USUALLY Averaging

(AKA. IDENTITY) $P=0.5$

WIDE MINIMA

COULD BE DUE TO OVERPARAM. (NETWORKS)

ENCOURAGES MULTIPLE MINIMA

ENCOURAGES MULTIPLE INTERNAL ENCOURAGINGS

INFEASIBLE

