

# OPTIMIZATION FOR TRAINING

**TRAINING PROBLEM:**  $J^*(\theta) = E_{(x,y) \sim P_{DATA}} L(\hat{f}(x;\theta), y)$  • ANN TRAINING IS DIFFERENT FROM TRADITIONAL OPTIMIZATION  $\rightarrow$  LOSS IS NOT THE TASK IN AND OF ITSELF

• WE DO NOT HAVE / KNOW TRUE DISTRIBUTION / DATA GENERATING PROCESS BUT A FINITE SET OF SAMPLES FROM IT  $\rightarrow$   
 $\rightarrow$  **EMPIRICAL RISK MINIMIZATION**  $E_{(x,y) \sim \hat{P}_{EMPIRICAL}} L(\hat{f}(x;\theta), y) = \frac{1}{M} \sum_{i=1}^M L(\hat{f}(x_i;\theta), y_i)$  AND HOPES TO IMPROVE ACTUAL RISK. PRONE TO OVERFITTING (MEMORIZING TRAINING SET)  
ALSO PROBLEMS BECAUSE MANY LOSSES, IE 0-1, HAVE NO USEFUL DERIVATIVES TO GUIDE GRADIENT DESCENT  
 $\rightarrow$  RARELY USED FOR DEEP LEARNING

$\rightarrow$  **SURROGATE LOSS FCNS** WHEN DERIVATIVES ARE COM OR TOO COMPLEX / EXPENSIVE TO CRUNCH. 0-1 LOSS  $\rightarrow$  LOG-LIKELIHOOD SURROGATE CAN EVEN LEARN MORE THAN ORIGINAL LOSS. ~~OR~~ ORIGINAL LOSS STILL USEFUL TO OBTAIN EARLY STOPPING. NO GUARANTEES THEY STOP IN LOCAL MINIMA BUT WHEN OVERFITTING STARTS INCREASING IN

$\rightarrow$  **BATCHES & MINIBATCHES** EXACT GRADIENT ON FULL TRAINING SET IS EXPENSIVE. STORER:  $\hat{G}/N \rightarrow$  LESS THAN LINEAR INCREASING IN USING MORE SAMPLES.  $\rightarrow$  USE SAMPLE SUBSET, USES DATA REDUNDANCY TO ITS ADVANTAGE. BATCH  $\rightarrow$  MINIBATCH  $\rightarrow$  ONLINE.  
• APPROPRIATE BATCH SIZING MAKES GOOD USE OF MULTICORE SETUPS.  $\rightarrow$  POWER OF 2 SIZE FOR GPU SETUPS. • SMALL BATCHES HAVE REGULARIZATION EFFECT.  
• **1<sup>ST</sup> ORDER METHODS:** BATCH SIZE  $\sim 100$  • **2<sup>ND</sup> ORDER METHODS:** LARGER,  $\sim 10,000$  ESPECIALLY IF  $H$  IS POORLY CONDITIONED.  
• SHUFFLE TRAINING SET TO DISCORRELATE SUCCESSIVE SAMPLES.  $\rightarrow$  MAKE MINIBATCHES MORE UNIFORM.  
• DIFFERENT MINIBATCH UPDATE CAN BE COMPUTED IN PARALLEL

$\rightarrow$  **GENERALIZATION ERROR:** IS UNBIASED (BUT NOISY) ESTIMATOR OF EXACT GENERALIZATION ERROR UNTIL END OF 1<sup>ST</sup> EPOCH. SO IS EQUIVALENT TO SGD. THEN THEY DIVERGE EXPLOITING FOR ONLINE-ONLY SETTINGS.

## ISSUES IN ANN OPTIMIZATION BECAUSE OBJECTIVE IS NONCONVEX

- **ILL-CONDITIONING** OF HESSIAN. 2<sup>ND</sup> ORDER UPDATE TERM DEPENDS ON EIGENVALS OF  $H$  AND ALIGNMENT OF EIGENVALS WITH  $g$ . WHILE  $g$  IS ALIGNED WITH LARGE POSITIVE EIGENVALS UPDATE MAY MOVE UPHILL. MONITOR WHERE  $g^T g$  AND  $g^T H g$ . 1<sup>ST</sup> STAYS SAME, 2<sup>ND</sup> GROWS FOR VERY SMALL LEARNING RATES  
• USUALLY NEWTON METHOD IS SUFFICIENT TO FIX THE ISSUE, NOT SO MUCH IN NNs  
• ALSO WHEN EIGENVALS OF OBJECTIVE FCN ARE CLOSE TO 0, ~~OR~~ INVERTING  $H$  BECOMES NUMERICALLY DIFFICULT AND AFFECTS NO OF EPOCHS FOR CONVERGENCE

- **LOCAL MINIMA** FITNESS LANDSCAPE OF NNs HAS MANY LOCAL MINIMA BECAUSE OF **MODEL IDENTIFIABILITY**. USUAL **WEIGHT-SPACE SYMMETRY** PROBLEM. ALSO WEIGHT SCALING SYMMETRY WHEN LINEAR OR MAXOUT UNITS (EVERY LM IS ON A MAXN HYPERPLANE OF EQUIVALENT LM,  $\alpha, 1/\alpha$ ) IF NO WEIGHT DECAY  
**HOWEVER** LM FROM UNIDENTIFIABILITY ALL HAVE SAME COST  $\rightarrow$  NOT A PROBLEM FOR GRADIENT DESCENT IF LM PROBLEMATIC IF HIGH-COST. NOWADAYS NOT BELIEVED TO BE COMMON. CHECK BY PLOTTING GRADIENT NORM OVER TIME

- **SADDLE POINTS & PLATEAUS** SADDLE  $\rightarrow$  H HAS BOTH POSITIVE AND NEGATIVE EIGENVALUES. EXPONENTIALLY MORE COMMON THAN LOCAL MINIMA IN HIGH-DIMENSIONAL SPACES !!  
IT HAPPENS THAT FOR RANDOM FCNS EVALS ARE MORE LIKELY TO BE POSITIVE IN LOW COST REGIONS. HIGH COST CRITICAL POINTS  $\rightarrow$  LIKELY A SADDLE  
• GRADIENT DESCENT SEEMS IMMUNE TO SADDLES • 2<sup>ND</sup> ORDER METHODS (SEEKING A JUMP TO CRITICAL POINT), NEWTON, CAN FALL IN THEM.  
• PROMISING 2<sup>ND</sup> ORDER SADDLE-FREE NEWTON  $\rightarrow$  STILL NOT SCALES • PLATEAUS ARE PROBLEMATIC  $\rightarrow$  IN NONCONVEX MAY BE HIGH-VALUE

- **EXPLODING GRADIENT:** COMMON IN RNN WITH LONG-TERM DEPENDENCIES. GRADIENT DESCENT WORKS WRT INFINITESIMAL MOVES. **GRADIENT CLIPPING!**  
- CLIFFS: ADJUST SIZE OF UPDATE STEP WHEN GRADIENT MAGNITUDE IS TOO HIGH. NEED THE DIRECTION.

- **ISSUES IN RNN W/ LONG-TERM DEPENDENCIES:** **EXPLODING / VANISHING GRADIENTS:** COMPOSITION OF GRADIENTS IS PRODUCT OF JACOBIANS THROUGH LAYERS/TIME. MULTIPLYING MANY LEADS TO V. LARGE / SMALL VALUES. STATISTICALLY UNLIKELY IF INDEPENDENT MATRICES. IN RNNs ARE RELATED  $\rightarrow$  SIMILAR EIGENVALUES, BLENDING OR CRASHES  
• FOR RNN TO LEARN, JACOBIANS MUST HAVE  $DET \leq 1$   $\rightarrow$  LEADS TO VANISHING! LONG TIME TO LEARN, IF AT ALL. MISSING BY SHORT TERM DEPENDENCIES, USUALLY 10-20 STEPS

- **INTRACTABLE GRADIENTS:** SUCH IT UP, USE APPROXIMATIONS AND PLAY

## BASIC OPTIMIZATION ALGORITHMS

- **STANDARD (BATCH) GRADIENT DESCENT:**  $\theta \leftarrow \theta + \epsilon \nabla \mathcal{L}$  LOOL, IMPRESSIVE CONVERGENCE. NO NEED TO REDUCE  $\epsilon$  OVER TIME IF TRUE BATCH. RARELY USED BECAUSE NOT EXPLOITS STRUCTURE OF PROBLEM
- **SGD** MOST USED IN ANN/DEEP LEARNING. TRUE ONLINE OR W/ MINIBATCHES.  $E[\hat{g}] = g$ . LEARNING RATE  $\eta$  MUST BE REDUCED IN EPOCHS. STATISTICALLY CONVERGES SLOWER THAN BATCH BUT GIVEN LIMITED COMPUTATIONAL RESOURCES IT INITIALLY CONVERGES MUCH FASTER. ALSO WE ALWAYS CARE ABOUT GEN. BOB SLOW WHEN GRADIENT IS SMALL
- **MOMENTUM** OPTIMAL WHEN GRADIENT IS CONSISTENT ON CONSECUTIVE MINIBATCHES  $v \leftarrow \alpha \cdot v + \eta \nabla \mathcal{L}, \theta \leftarrow \theta + v$
- **NESTEROV MOMENTUM:** GRADIENT EVALUATED **AFTER** APPLYING CURRENT VELOCITY. MAKES CONVERGENCE FASTER IN BATCH CASE

## ADAPTIVE LEARNING RATE ALGORITHMS

LEARNING RATE IS DOMINANT FACTOR FOR CONVERGENCE.

- **DELTA-RULE-DELTA** FOR FULL BATCH ONLY. PER-PARAMETER LOSS DERIVATIVE. IF SAME INCREASE LR. IF DIFF SIGN DECREASE SIGN
- **ADAGRAD** PER PARAM LR SCALED INVERSELY TO SUM OF SQUARED PARTIAL DERIVATIVES OVER ITERATIONS  $\rightarrow$  UNDER PD, RAPID DECREASE IN LR. MORE PROGRESS WHERE SLOPE IS UPSTAIR. FOR DEEP NETS ACCUMULATION OF SQUARED GRADIENTS FROM SAME RESULTS IN BAD EFFECTIVE LR, TOO SLOW.
- **RMSPROP** ADAGRAD BUT GRADIENT ACCUMULATION CHANGE INTO EXPONENTIALLY WEIGHTED MOVING AVERAGE WITH WINDOW  $\rho$  BECAUSE RELATIVE SLOPE OF DERIVATIVES MIGHT CHANGE AS TRAINING PROGRESSES. CAN BE COMBINED W/ NESTEROV. EASY. CURRENTLY IN FASHION.  $R \leftarrow \rho R + (1-\rho)g^2$
- **ADAM** SIMILAR TO RMSPROP + MOMENTUM BUT MOMENTUM IS ESTIMATE OF 1ST ORDER MOMENT OF GRADIENT (EMA). ALSO CORRECTIONS TO 1ST AND 2ND ORDER MOMENTS TO ACCOUNT FOR ORIGIN INITIALIZATION.  $s \leftarrow \rho_1 s + (1-\rho_1)g$ ;  $R \leftarrow \rho_2 R + (1-\rho_2)g^2$ ;  $\hat{s} \leftarrow s / (1-\rho_1^t)$ ;  $\hat{R} \leftarrow R / (1-\rho_2^t)$ ;  $\Delta \theta = -\alpha \hat{s} / \sqrt{\hat{R} + \epsilon}$
- **ADADelta** ALSO TRIES TO FIX ISSUES WITH ADAGRAD. INCORPORATES 2ND ORDER INFORMATION. SQUARE ROOTS OF EMA OF INCREMENTAL PARTIAL DERIVATIVES SQUARES. 
$$\Delta \theta = -\frac{\sqrt{s + g^2}}{\sqrt{R + g^2}} \cdot g$$
  $s \leftarrow \rho s + (1-\rho)[\Delta \theta]^2$
- **NO CLEAR WINNER!** RMSPROP + ADADelta ARE REQUEST BUT SGD (MOM), RMSPROP (MOM), ADADelta, ADAM ARE ALL GOOD.

## APPROXIMATE 2ND ORDER METHODS

FOR EMPIRICAL RISK BUT ALSO WORK FOR OBJECTIVES WITH REGULARIZATION TERMS

- **NEWTON'S METHOD**  $\rightarrow$  2ND ORDER TAYLOR APPROXIMATION  $\theta^* = \theta_0 - [H(\theta_0)]^{-1} \nabla \ell(\theta_0)$ . JUMPS DIRECTLY TO MINIMA FOR LOCALLY QUADRATIC FCN. ELSE FIXED POINT UPDATE. (AS LONG AS H POSITIVE DEFINITE) EM-LINE (INVERTS H  $\rightarrow$  UPDATE  $\theta$ ) • **IN DEEP LEARNING!** SADDLE PROBLEM  $\rightarrow$  REGULARIZE IT  $\theta^* = \theta_0 - [H(\theta_0) + \lambda I]^{-1} \nabla \ell(\theta_0) \rightarrow$  **LEVENBERG-MARQUARDT** OR AS LONG AS NEGATIVE EIGENVALS ARE CLOSE TO 0, ELSE NUMBERS MAKE IT SAME AS 1ST ORDER
- **COMPUTATIONALLY VERY EXPENSIVE**  $O(n^3)$ !!!
- **CONJUGATE GRADIENTS** AVOIDS HESSIAN INVERSION. NORMAL STEEPEST DESCENT HAS ORTHOGONAL SUCCESSIVE STEPS  $\rightarrow$  ZIGZAGGING.  $\rightarrow$  UNDOES PROGRESS' MINIMUM ALONG PREVIOUS DIRECTIONS IS NOT PRESERVED. **CG FIXES THIS**  $d_t = \nabla \ell(\theta) + \beta d_{t-1}$ , ADDS BACK SOME OF PREVIOUS SEARCH DIRECTION.
  - $d_t^T H(\theta) d_{t-1} = 0 \rightarrow$  IS STEEPEST DESCENT IN  $\phi$  SPACE ( $\phi = A^{1/2} Q^T \theta$ ) 'HESSIAN' SPACE [ $H = Q \Lambda Q^T$  ORTHOGONAL  $Q$ ]
  - CAN COMPUTE  $\beta$  WITHOUT EIGENDECOMPOSITION OF  $H$  OR  $H$  AT ALL: FLETCHER-REEVES, POWELL-REINHOLD FORMULAS ONLY BECOM ON)
  - $\rightarrow$  ONLY REQUIRE  $K$  LINE SEARCHES (ITERATIONS) IN  $n$ -DIM SPACE **NONLINEAR CG:** OCCASIONALLY DOES A STANDARD LINE SEARCH STEP IN DL: START WITH A FEW STANDARD SGD ITERATIONS, THEN DO CONJ. GRADS
- **BFGS** APPROXIMATES  $H^{-1}$  WITH LOW RANK APPROXIMATION  $M_t$  REFINED ITERATIVELY  $M_t = M_{t-1} + \dots$ . RANK ONE UPDATES  $\theta^{R*}$ : UPDATE IS  $O(n^2)$   $\theta_t = M_t g_t$  LINE SEARCH FOR STEP SIZE  $\eta^* = \text{ARGMIN}_\eta (\ell(\theta_t + \eta g_t))$ ,  $\theta_{t+1} = \eta^* g_t$ . LESS TIME TO REFINED LINE SEARCH STEP BUT STORAGE OF  $M$  IS IMPRACTICALLY LARGE FOR **TOUGH STUFF**
- **L-BFGS** REDUCES  $M_{t-1}$  WITH  $l$ , UPDATES  $\theta_t = -g_t + b \Delta + \alpha \phi$ .  $\Delta, \phi = g_t - g_{t-l}$ ,  $\alpha, b$  DEPEND ALSO ON THEM. IF EXACT LINE SEARCHES: MUTUALLY CONJUGATE DIRECTIONS BUT ALSO ON IF APPROX LINE SEARCH



# NATURAL GRADIENT METHODS

ATTEMPT TO MAKE OPTIMIZATION INVARIANT TO MODEL PARAMETRIZATION  $\theta$ . STEEPEST DIRECTION IN SPACE OF PROBABILITY DISTRIBUTIONS OUTPUT BY MODEL.

$$\Delta_{NL} = \underset{\Delta}{\text{ARGMIN}} E_{\text{DATA}} [-\log p_{\theta+\Delta}(x)] = KL(p_{\theta}(x) || p_{\theta+\Delta}(x)) = \Delta KL \quad KL(p_{\theta} || p_{\theta+\Delta}) = -\frac{1}{2} \Delta^T E_{\text{FO}} [\nabla^2 \log p_{\theta}] \Delta \quad | \quad -\nabla^2 \log p_{\theta} = \text{FISHER INFO MATRIX}$$

$$L_N(\theta, \Delta\theta) = E_{\text{DATA}} [-\log p_{\theta}] + E_{\text{DATA}} [-\nabla \log p_{\theta}]^T + \frac{\lambda}{2} \Delta^T E_{\text{FO}} [-\nabla^2 \log p_{\theta}] \Delta \quad \theta_{t+1} = \theta_t + (E_{\text{FO}} [-\nabla^2 \log p_{\theta}])^{-1} E_{\text{DATA}} [-\nabla \log p_{\theta}] \quad \bullet \text{ SCALES WITH INVERSE OF FISHER INFO MATRIX}$$

EXPECTATIONS AVERAGE OVER MODEL DISTRIBUTIONS, NEWTON'S DOES IT OVER DATASET

## OPTIMIZATION STRATEGIES

**COORDINATE DESCENT:** MINIMIZE WRT ONE VAR AT TIME, OR GROUPS (BLOCK COORD. DESCENT). ON WHEN VARS/GROUPS RELATIVELY ISOLATED (THE SPARSE CODING DICTIONARY AND CODES)

**INITIALIZATION:** BL IS VERY SENSITIVE TO IT, STILL SIMPLY HEURISTIC BECAUSE WE DON'T HAVE ANY BETTER, ALSO OPTIMIZATION VS GENERALIZATION TRADEOFF

- **SYMMETRY BREAKING!** SAME FCN ON SAME INPUTS NEED DIFFERENT INIT PARAMS ELSE THEY'LL EVOLVE THE SAME. NO USE + FWD/BWD NULLSPACES  
→ MOTIVATES RANDOM INITIALIZATION FROM GAUSSIAN OR UNIFORM. BIAS ARE HEURISTIC CONSTANT.
- IDEALLY LARGE VALUES HELP IN BREAKING SYMMETRY BUT LEAD TO OTHER ISSUE: GRADIENT EXPLOSION, CHAOS, SATURATION
- CONSIDERATION OF OPTIMIZATION VS GENERALIZATION. WE ALSO WANT SMALL PARAMS B/C OF REGULARIZATION.
- **INIT IS EFFECTIVELY AN IMPLICIT REG.** WE IMPOSE ON THE MODEL. (CLOSE TO 0, LARGE VALUES...)
- SAMPLING HEURISTICS WRT M INPUTS, N OUTPUTS, ORTHOGONAL MATRICES, SCALE INIT. BY CONSTANT GAIN FACTOR  $g$  (HIGHWAY!); PRESERVING NORM.  
• STANDARD VALUES ARE 1 OVER MANY STEPS
- **SPARSE INITIALIZATION** (K-ARM ZERO WEIGHTS) - NOT GOOD ON MAXOUT UNITS
- → FRAME IT AS AN HYPERPARAMETER • MANUALLY ADJUST BY LOOKING AT STDDEV OR RANGE OF ACTIVATIONS LAYER-BY-LAYER SCALE
- **BIAS INITIALIZATION:** DEPENDS ON WEIGHT. COMMONLY 0. IF CREATIVE MODEL → MATCH MARGINALS • RELU → SMALL POSITIVE BIAS TO AVOID SATURATION OUTPUTS
- LSTM GATES → 1 TO HAVE THEM INITIALLY OFF
- **VARIANCE/PRECISION:** SET TO 1 OR TO MAXIMAL VARIANCE
- USE OTHER ML METHODS / SUPERVISED TRAINING

**GREEDY SUPERVISED PRE-TRAINING:** START BY TRAINING A SIMPLER MODEL, OR ATTEMPT A SIMPLER TASK • **GREEDY ALGOS + FINE-TUNING** ON JOINT PROBLEM  
IT GIVES GUIDANCE TO MIDDLE LAYERS!  
• ONE LAYER AT A TIME USING PREV LAYER OUTPUTS AS INPUTS • USE FIRST AM. BEST TO INIT MIDDLE OF AN EVEN DEEPER LAYER  
• TRANSFER LEARNING; TRAIN ON A TASK, ADD LAYERS AND TRY ON ANOTHER TASK

**DESIGNING MODELS:** MORE IMPORTANT / BETTER TO CHOOSE A MODEL EASY TO OPTIMIZE THAN A POWERFUL OPTIMIZATION ALGO  
HISTORICAL IDEAS: FORWARD NETS WITH LINEAR UNITS AND DIFFERENTIABLE ACTIVATIONS BECAUSE OPTIMIZATION GETS EASIER.  
• LINEAR PATHS, SKIP CONNECTIONS, OR MULTIPLE TRAINING HEADS TO BOOST GRADIENT SIGNAL TO DISCARD IN APPLICATION.

**CONTINUATION METHODS:** SEQUENCES OF OBJECTIVE FUNCTIONS IN THE SAME PARAMETERS, INCREASINGLY DIFFICULT. (AND WELL BEHAVING ON  $\theta$  SPACE.) ORIGINALLY DESIGNED TO HELP AVOID LOCAL MINIMA (BESIDES MIN). **BLURRING** → CONVOLVING WITH A GAUSSIAN, POLYMERIZATION TEMPERATURE SCHEME. HOPE IS MAKING THE FCN MORE CONVEX. **CURRICULUM LEARNING!** LEARN BASIC TASKS, THEN MOVE TO MORE COMPLEX TASKS REQUIRING MASTERY OF THE BASICS  
SUCCESSFUL IN MLP AND VISION TASKS