# ML TUTORIAL 2 — BACKPROPAGATION, OPTIMIZERS

HOW TO LEARN WEIGHTS OF AN MLP?

$$\bar{y} = \dots \sigma\left(\sigma\left(\bar{x}\,\underline{W_1} + \bar{b}_1\right)\underline{W_2} + \bar{b}_2\right)\dots$$

$$E = \frac{1}{2}\left(\bar{y} - \hat{\bar{y}}\right)^T\left(\bar{y} - \hat{\bar{y}}\right)$$

<span style="color:red">← MSE LOSS, VECTORIZED</span>

WE NEED $\dfrac{\partial E}{\partial \underline{W_1}}$, $\dfrac{\partial E}{\partial \underline{W_2}}$, $\dots$, $\dfrac{\partial E}{\partial \bar{b}_1}$, $\dfrac{\partial E}{\partial \bar{b}_2}$ $\dots$,

THESE GRADIENTS ARE THEN USED FOR GRADIENT DESCENT, TO MODIFY THE CURRENT PARAMETERS:

$$\underline{W_1}^{NEW} = \underline{W_1}^{OLD} - \eta\,\frac{\partial E}{\partial \underline{W_1}}$$

## IDEA: USE CHAIN RULE

EXAPLE:

$$y = h\left(g\left(f(x)\right)\right)$$

IMMEDIATE VARIABLES

$$u = f(x)$$
$$v = g\left(f(x)\right)$$

$$\frac{dy}{dx} = \frac{dy}{dv}\,\frac{dv}{du}\,\frac{du}{dx}$$

ANOTHER NOTATION THAT CAN BE USED IN CASE IT IS CLEAR FROM THE CONTEXT WHICH VARIABLE TO USE:

$$\frac{d\,h\left(g\left(f(x)\right)\right)}{dx} = h'\left(g\left(f(x)\right)\right)\,g'\left(f(x)\right)\,f'(x)$$

WHERE $f'(x) = \dfrac{df(x)}{dx}$ , $g'(x) = \dfrac{dg(x)}{dx}$ , ...

FROM THE NOTATION

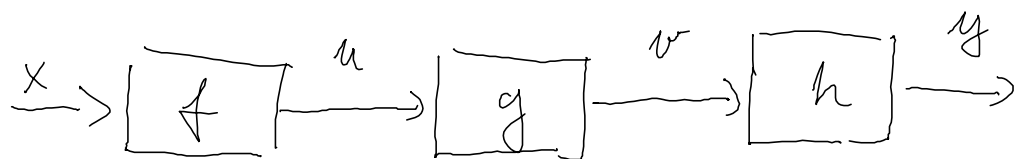$$h'(\overbrace{g(f(x))}^{v}) \; g'(\overbrace{f(x)}^{u}) \; f'(x)$$

IT CAN BE SEEN, THAT THE PROCESS CAN BE DIVIDED IN 2 PARTS:

①  CALCULATE THE VALUE OF ALL IMMEDIATE VARIABLES AND SAVE THEM:

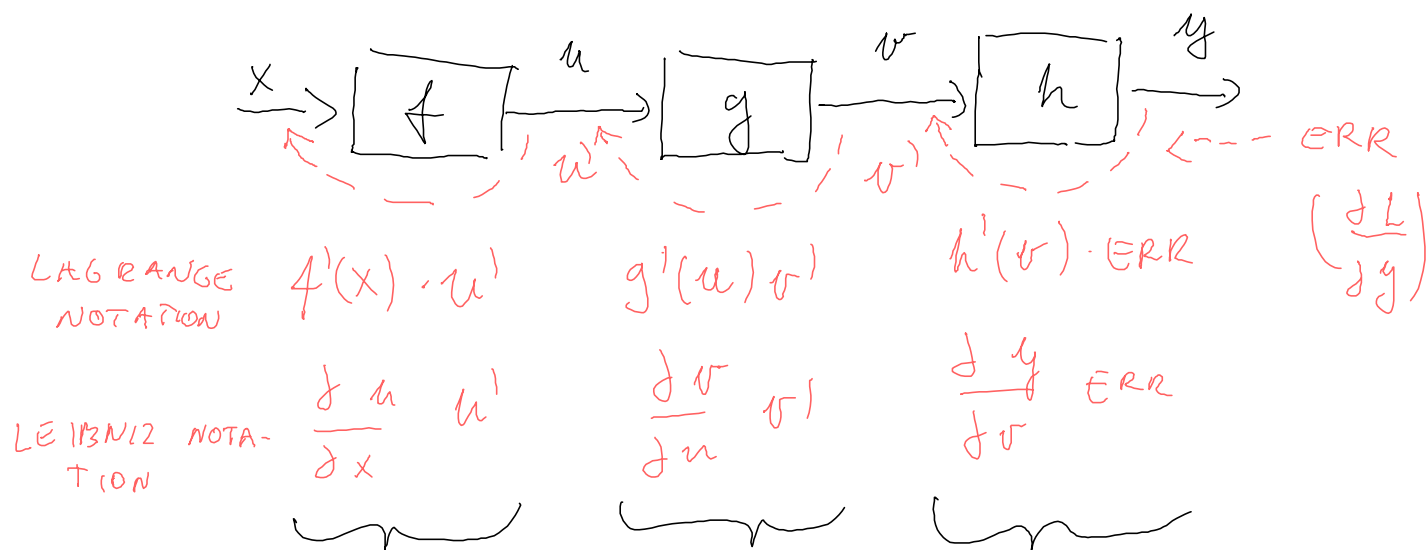$$u = f(x) \qquad v = g(u) \qquad y = h(v)$$

THIS IS THE  FORWARD  PASS

COMPUTATION GRAPH:

$$x \longrightarrow \boxed{f} \xrightarrow{\;u\;} \boxed{g} \xrightarrow{\;v\;} \boxed{h} \xrightarrow{\;y\;}$$

(2) CALCULATE THE DERIVATIVE OF EACH FUNCTION INDEPENDENTLY AT THE POINT OF THEIR ARGUMENT IN FORWARD PASS. THEN MULTIPLY THEM TOGETHER.

THIS IS THE BACKWARD PASS



LAGRANGE NOTATION: $f'(x) \cdot u'$   $g'(u) v'$   $h'(v) \cdot ERR$   $\left( \dfrac{dL}{dy} \right)$

LEIBNIZ NOTATION: $\dfrac{du}{dx}\ u'$   $\dfrac{dv}{du}\ v'$   $\dfrac{dy}{dv}\ ERR$

NOTE: ALL OF THEM ARE LOCAL COMPUTATION. THEY ONLY DEPEND ON THE ARGUMENT OF THE FORWARD PASS, AND THE ERROR TERM OF THE BACKWARD PASS

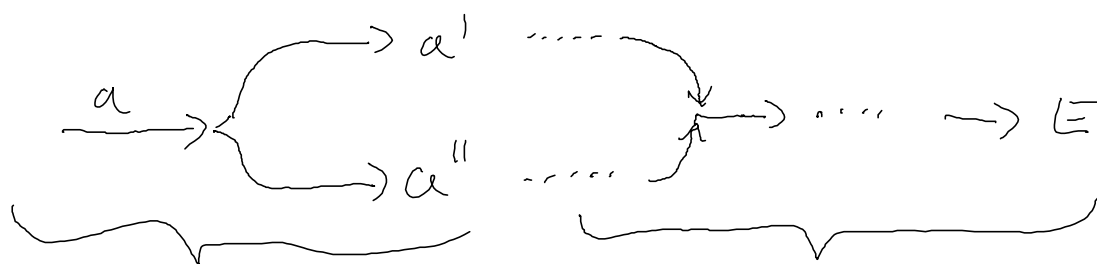ARGUMENTS (ACTIVATIONS) OF FORWARD PASS  $\longrightarrow$

GRADIENTS FROM THE BACKWARD PASS  $\longleftarrow$

THE LOCAL COMPUTATION DOESN'T CARE HOW
THE INPUT IS CREATED OR HOW THE OUTPUT
GRAD IS FORMED; THEY ALWAYS HAVE THE SAME
FORM. THIS ENABLES THE FUNCTIONS TO
BE CHAINED TOGETHER LIKE BUILDING BLOCKS,
INDEPENDENTLY OF EACH OTHER. PYTORCH / TENSOR-
FLOW IS JUST A LIBRARY OF SUCH BLOCKS.
IN PRACTICE, WE USE A VECTORIZED FORM.

RECIPE:

① START FROM THE BACK, GO TO THE FRONT

② CALCULATE GRADS FOR ALL INPUTS OF THE BLOCK

③ IF THERE IS A FORK, ADD GRADS TOGETHER (MULTI-VARIATE CHAIN RULE)



FORK
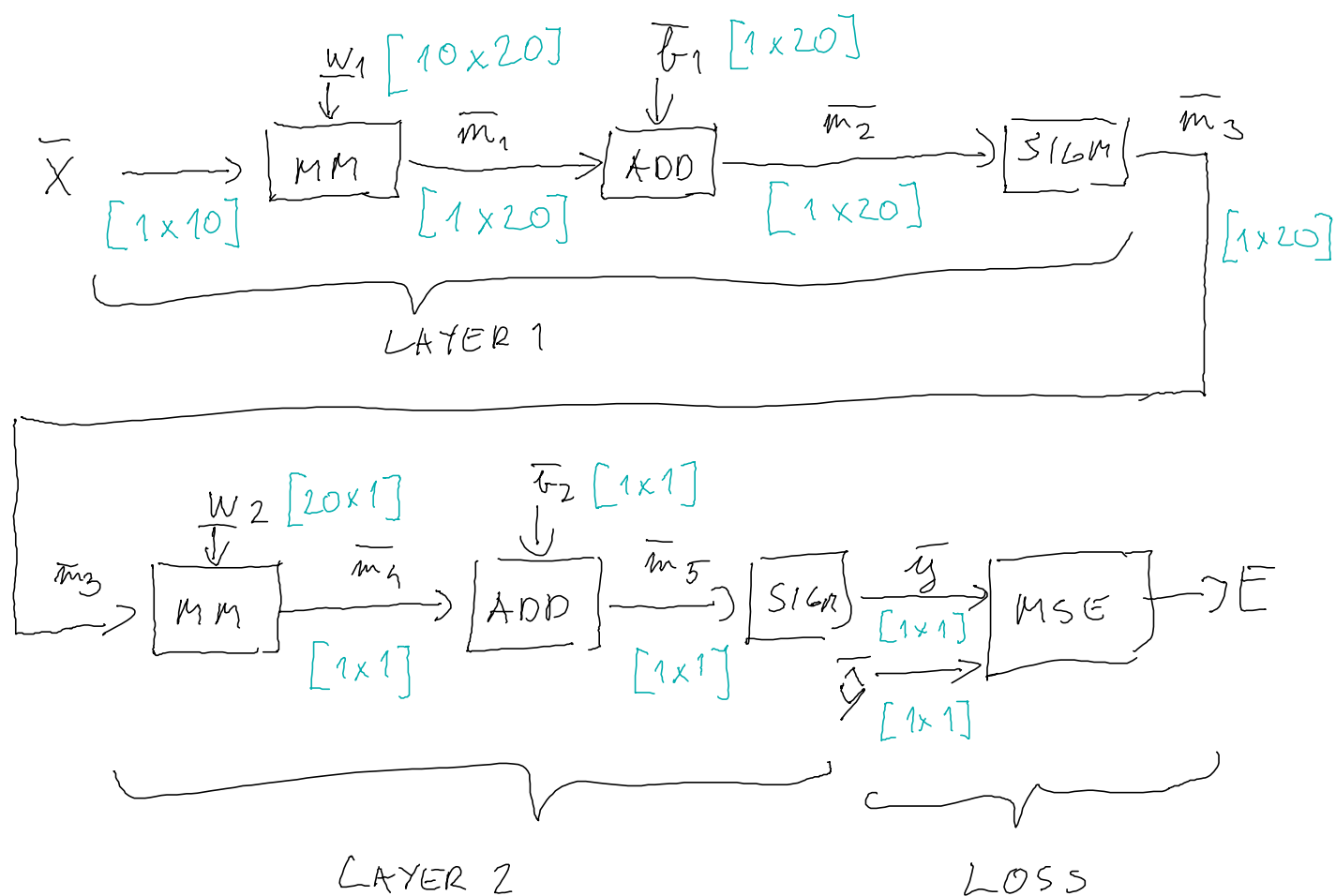
$$\frac{\partial E}{\partial a} = \frac{\partial E}{\partial a'} + \frac{\partial E}{\partial a''}$$

THIS PART ALWAYS EXISTS, BECAUSE WE HAVE A SINGLE, SCALAR "E".

NOTE: YOU WILL ENCOUNTER ALL THE VARIABLES YOU ARE INTERESTED IN DURING THE BACKWARD PASS. THEY WILL BE ON DIFFERENT BRANCHES. BUT THEIR "END" WILL BE THE SAME AND HAVE TO BE COMPUTED ONLY ONCE.

# LONG EXAMPLE: 2 LAYER MLP

$\overline{X}$ $[1 \times 10]$ $\rightarrow$ [MM] $\overline{m_1}$ $[1 \times 20]$ $\rightarrow$ [ADD] $\overline{m_2}$ $[1 \times 20]$ $\rightarrow$ [SIGM] $\overline{m_3}$ $[1 \times 20]$

$W_1 [10 \times 20]$ (into MM)

$\overline{b_1} [1 \times 20]$ (into ADD)

LAYER 1

$\overline{m_3}$ $\rightarrow$ [MM] $\overline{m_4}$ $[1 \times 1]$ $\rightarrow$ [ADD] $\overline{m_5}$ $[1 \times 1]$ $\rightarrow$ [SIGM] $\overline{y}$ $[1 \times 1]$ $\rightarrow$ [MSE] $\rightarrow E$

$W_2 [20 \times 1]$ (into MM)

$\overline{b_2} [1 \times 1]$ (into ADD)

$\overline{d}$ $[1 \times 1]$ (into MSE)

LAYER 2        LOSS

WE WANT TO CALCULATE:

$$\frac{\partial E}{\partial W_1}, \quad \frac{\partial E}{\partial \overline{b_1}}, \quad \frac{\partial E}{\partial W_2}, \quad \frac{\partial E}{\partial b_2}$$

WHAT IS THE PATH WE HAVE TO TAKE?

$\dfrac{\partial E}{\partial W_1}$ :  $E \rightarrow \overline{y} \rightarrow \overline{m_5} \rightarrow \overline{m_4} \rightarrow \overline{m_3} \rightarrow \overline{m_2} \rightarrow \overline{m_1} \rightarrow W_1$

$\dfrac{\partial E}{\partial \overline{b_1}}$ :  $E \rightarrow \overline{y} \rightarrow \overline{m_5} \rightarrow \overline{m_4} \rightarrow \overline{m_3} \rightarrow \overline{m_2} \rightarrow \overline{b_1}$

SAME: CALCULATE ONCE — DYNAMIC PROGRAMMING

START FROM THE END, MODULE BY MODULE:

$y \rightarrow$ $\boxed{MSE}$ $\rightarrow E$   FORWARD:

$\tilde{y} \rightarrow$

$$E = \frac{1}{2}(y - \tilde{y})^2$$

BACKWARD:

$$\frac{\partial E}{\partial y} = \frac{\partial}{\partial y} \frac{1}{2}(y - \tilde{y})^2 = y - \tilde{y}$$

RESULT SOFAR:

$$\frac{\partial E}{\partial y} = y - \tilde{y}$$

- - - - - - - - - - - -

$\overline{m_5} \rightarrow \boxed{SIGM} \rightarrow \tilde{y}$   FORWARD:

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

BACKWARD:

$$\frac{\partial}{\partial x}(1 + e^{-x})^{-1} = -(1 + e^{-x})^{-2} \frac{\partial}{\partial x}(1 + e^{-x}) =$$

$\uparrow$
CHAIN RULE

$$= -(1 + e^{-x})^{-2}\left(0 + \frac{\partial}{\partial x} e^{-x}\right) =$$

$$= -(1 + e^{-x})^{-2}(-1)e^{-x} =$$

CHAIN RULE AGAIN

$$= (1 + e^{-x})^{-2} e^{-x} = \frac{e^{-x}}{(1 + e^{-x})^2} =$$

$$= \frac{1}{1 + e^{-x}} \cdot \frac{e^{-x} \boxed{+ 1 - 1}}{1 + e^{-x}} =$$

$$= \frac{1}{1 + e^{-x}} \cdot \frac{e^{-x} + 1 - 1}{1 + e^{-x}} =$$

$$= \frac{1}{1 + e^{-x}} \cdot \left( \frac{1 + e^{-x}}{1 + e^{-x}} - \frac{1}{1 + e^{-x}} \right) =$$

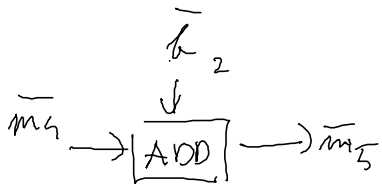$$= \frac{1}{1 + e^{-x}} \cdot \left( 1 - \frac{1}{1 + e^{-x}} \right) =$$

$$= \boxed{\sigma(x) \left( 1 - \sigma(x) \right)}$$

RESULT  SOFAR

$$\frac{\partial C}{\partial \overline{m_5}} = \frac{\partial E}{\partial \overline{y}} \cdot \frac{\partial \overline{y}}{\partial \overline{m_5}}$$

ORDER OF MULTIPLICATION DOESN'T MATTER. IT IS ELEMENTWISE

$$\frac{\partial E}{\partial \overline{m_5}} = \sigma(\overline{m_5}) \left( 1 - \sigma(\overline{m_5}) \right) \left( y - \hat{y} \right)$$

$$\overline{m_4} \longrightarrow \boxed{ADD} \longrightarrow \overline{m_5}$$

$$\overset{\overline{b}_2}{\downarrow}$$

FORWARD:
$$\overline{m_5} = \overline{m_4} + \overline{b}_2$$

BACKWARD:

2 OUTPUTS:

$$\frac{\partial \overline{m_5}}{\partial \overline{m_4}} = 1 \qquad \frac{\partial \overline{m_5}}{\partial \overline{b}_2} = 1$$

RESULTS SO FAR:

$$\frac{\partial E}{\partial \overline{b}_1} = \frac{\partial E}{\partial \overline{m_5}} \cdot \frac{\partial \overline{m_5}}{\partial \overline{b}_1} = \frac{\partial E}{\partial m_5}$$

$$\boxed{\frac{\partial E}{\partial \overline{b}_1} = \sigma(\overline{m_5})\left(1 - \sigma(\overline{m_5})\right)\left(y - \hat{y}\right)}$$

$\uparrow$ THIS IS ALREADY ONE OF THE GRADIENTS WE ARE INTERESTED IN.

$$\frac{\partial E}{\partial \overline{m_4}} = \frac{\partial E}{\partial \overline{m_5}} \cdot \frac{\partial \overline{m_5}}{\partial \overline{m_4}} = \frac{\partial E}{\partial \overline{m_5}}$$

$$\frac{\partial E}{\partial \overline{m_4}} = \sigma(\overline{m_5})\left(1 - \sigma(\overline{m_5})\right)\left(y - \hat{y}\right)$$
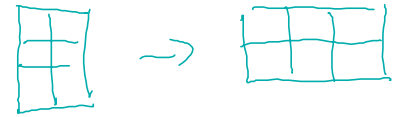
FORWARD

$$\overline{m}_4 = \overline{m}_3 \times \underline{W}_2$$

BACKWARD: 2 BRANCHES

$$\frac{\partial \overline{m}_4}{\partial \underline{W}_2} = \overline{m}_3 \; \textcircled{T}$$

TRANSPOSE:

$$\frac{\partial \overline{m}_4}{\partial \overline{m}_3} = \underline{W}_2^T$$

CHAIN RULE IS NONTRIVIAL. ORDER MATTERS

$$\frac{\partial E}{\partial \underline{W}_2} = \overline{m}_3^T \times \frac{\partial E}{\partial \overline{m}_4}$$

MATRIX MULTIPLICATION
ORDER MATTERS

$$\frac{\partial E}{\partial \overline{m}_3} = \frac{\partial E}{\partial \overline{m}_4} \times \underline{W}_2^T$$

np.matmul $\left( \frac{\partial Y}{\partial \overline{m}_4} , W2.T \right)$

A TRICK TO FIGURE OUT THE RIGHT ORDER
IS TO WRITE DOWN AN EXAMPLE WITH A
NON-SQUARE MATRIX, AND ENSURE THAT
THE SHAPE OF GRADIENT MATCHES THE
SHAPE OF THE INPUT.

NOTE 2: ALL OTHER OPERATIONS IN THIS
EXAMPLE ARE ELEMENTWISE. EACH
COMPONENT OF THE ACTIVATION VECTOR
CAN BE COMPUTED SEPARATELY, AS
THEY WERE SCALARS. THIS IS NOT TRUE
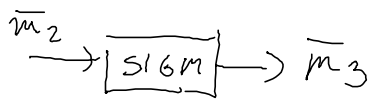FOR MATRIX MULTIPLICATION.

RESULTS SO FAR:

$$\frac{\partial E}{\partial \underline{W}_2} = \bar{m}_3^T \times \frac{\partial E}{\partial \bar{m}_5}$$

$$\frac{\partial E}{\partial \underline{w}_2} = \bar{m}_3^T \times \left( \sigma(\bar{m}_5) \left( 1 - \sigma(\bar{m}_5) \right) \left( y - \hat{y} \right) \right)$$

THIS IS ONE OF THE GRADS WE
ARE INTERESTED IN

$$\frac{\partial E}{\partial \overline{m}_3} = \frac{\partial E}{\overline{m}_4} \times \underline{W}_2^T$$

$$\frac{\partial E}{\partial \overline{m}_3} = \left[ \sigma(\overline{m}_5)\left(1 - \sigma(\overline{m}_5)\right)\left(y - \hat{y}\right) \right] \times \underline{W}_2^T$$

---

$$\overline{m}_2 \longrightarrow \boxed{SIGM} \longrightarrow \overline{m}_3$$

FROM BEFORE:

$$\frac{\partial \overline{m}_3}{\partial \overline{m}_2} = \sigma(\overline{m}_2)\left(1 - \sigma(\overline{m}_2)\right)$$

<u>RESULT SO FAR</u>

$$\frac{\partial E}{\partial \overline{m}_2} = \frac{\partial E}{\partial \overline{m}_3} \cdot \frac{\partial \overline{m}_3}{\partial \overline{m}_2}$$

$$\frac{\partial E}{\partial \overline{m}_2} = \left[ \left[\sigma(\overline{m}_5)\left(1 - \sigma(\overline{m}_5)\right)\left(y - \hat{y}\right)\right] \times \underline{W}_2^T \right] \cdot$$

$$\cdot \ \sigma(\overline{m}_2)\left(1 - \sigma(\overline{m}_2)\right)$$

$$\overline{m}_1 \xrightarrow{\quad} \boxed{\text{ADD}} \xrightarrow{\overline{m}_2} \qquad \overline{b}_1 \downarrow$$

$$\frac{\partial \overline{m}_2}{\partial \overline{b}_1} = 1 \qquad\qquad \frac{\partial \overline{m}_2}{\partial \overline{m}_1} = 1$$

## RESULTS SO FAR

$$\frac{\partial E}{\partial \overline{b}_1} = \frac{\partial E}{\partial \overline{m}_2} \cdot \frac{\partial \overline{m}_2}{\partial \overline{b}_1}$$

$$\boxed{\frac{\partial E}{\partial b_2} = \left[ \left[ \sigma(\overline{m}_5)(1 - \sigma(\overline{m}_5))(y - \hat{y}) \right] \times \underline{w}_2^T \right] \cdot \, \sigma(\overline{m}_2)\left(1 - \sigma(\overline{m}_2)\right)}$$

ONE OF THE GRADIENTS WE ARE INTERESTED IN

$$\frac{\partial E}{\partial \overline{m}_1} = \frac{\partial E}{\partial \overline{m}_2} \cdot \frac{\partial \overline{m}_2}{\partial \overline{b}_1}$$

$$\frac{\partial E}{\partial \overline{m}_1} = \left[ \left[ \sigma(\overline{m}_5)(1 - \sigma(\overline{m}_5))(y - \hat{y}) \right] \times \underline{w}_2^T \right] \cdot \, \sigma(\overline{m}_2)\left(1 - \sigma(\overline{m}_2)\right)$$

$$\underline{X} \rightarrow \boxed{MM} \xrightarrow{\downarrow W_1} \overline{m}_1$$

<u>FROM BEFORE</u>

$$\frac{\partial \overline{m}_1}{\partial \underline{W}_1} = \overline{X}^T$$

WE CAN IGNORE $\dfrac{\partial \overline{m}_1}{\partial \overline{X}}$, WE DON'T NEED GRAD FOR THE INPUTS

$$\frac{\partial E}{\partial \underline{W}_1} = \widetilde{X}^T \times \frac{\partial E}{\partial m_1}$$

$$\frac{\partial E}{\partial \underline{W}_1} = \overline{X}^T \times \left[ \left[ \left[ \sigma(\overline{m}_5)(1 - \sigma(\overline{m}_5))(y - \hat{y}) \right] \times \underline{W}_2^T \right] \cdot \sigma(\overline{m}_2)(1 - \sigma(\overline{m}_2)) \right].$$

## <u>NOTES</u>

- THE GRADIENTS IN PRACTICE ARE NOT EXPRESSED SYMBOLICALLY AS WE DID HERE. THAT WOULD REQUIRE EACH OF THE EXP-RESSIONS TO BE CALCULATED SEPARATELY - <u>SLOW</u>. INSTEAD THE <u>VALUE</u> OF THE GRADIENTS IS CALCULATED AT EACH MODULE AND PASSED TO THE NEXT. THIS WAY EACH GRAD

IS CALCULATED ONLY ONCE. THIS IS EFFICIENT.

## DIFFERENCE BETWEEN GD AND BACKPROP

BACKPROPAGATION IS AN EFFICIENT DYNAMIC
PROGRAMMING ALGORITHM FOR CALCULATING
THE GRADIENT OF PARAMETERS WITH
RESPECT TO AN OUTPUT LOSS.

GRADIENT DESCENT IS AN ALGORITHM FOR
USING GRADIENTS (E.G. THE OUTPUTS OF
BACKPROP) TO IMPROVE THE MODELS
PERFORMANCE (DECREASING THE LOSS).
IT IS A WAY TO CHANGE THE MODELS
PARAMETERS BASED ON THEIR GRADIENTS,

# TRAINING TRICKS

## MOMENTUM

SGD IS VERY SLOW. AN EASY WAY TO IMPROVE IT IS THE MOMENTUM.

MOMENTUM OPTIMIZER ACCUMULATES THE GRADIENT AND DOES THE PARAMETER UPDATE WITH THE ACCUMULATED GRADS. THIS IS LIKE MOMENTUM IN PHYSICS.

WEIGHTS ARE INDEPENDENT IN PERSPECTIVE OF THE OPTIMIZER. $i$-th WEIGHT IS $w_i$

$$v_i^{NEW} = \mu \, v_i^{OLD} + \frac{\partial E}{\partial w_i}$$

$$w_i^{NEW} = w_i - \eta \cdot \boxed{\overline{v_i}} \quad - \text{ IN SGD THIS TERM IS } \frac{\partial E}{\partial w_i}$$
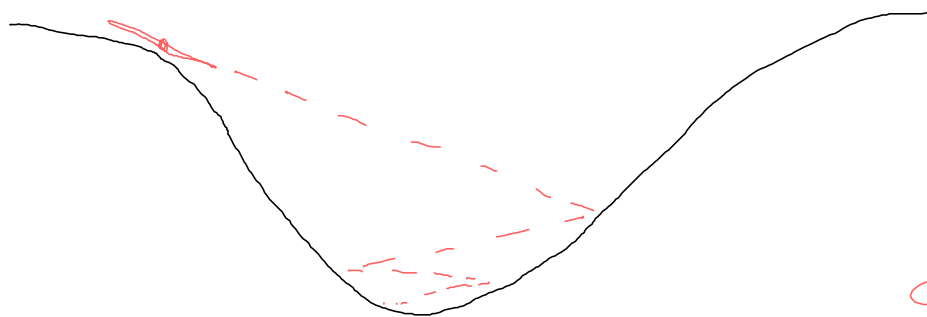
LEARNING RATE (TYPICALLY $0.1 \ldots 0.0001$)

$\mu$ - MOMENTUM (TYPICALLY $0.9 \ldots 0.99$)

INTUITION: IF GRADIENT CONSISTENTLY POINTS TO A DIRECTION, WE CAN ACCELERATE IN THAT DIRECTION (THE DIRECTION IS UNLIKELY TO CHANGE). HOWEVER IF THE GRADIENT CHANGES DIRECTION FREQUENTLY, WE SHOULD GO SLOWER IN ORDER NOT TO OSCILLATE NEAR THE MINIMUM.

SIGN IS CONSISTENT, ACCELERATE!

THE GRAD CHANGES DIRECTION SLOW DOWN!

# WEIGHT DECAY

COMMONLY USED REGULARIZER. PENALIZES THE MAGNITUDE OF THE WEIGHTS.

INTUITION: MAKE THE UNNECESSARY WEIGHTS 0.

(IF $\mathbb{E}[GRAD] = 0$, THERE IS NO DRIVING FORCE FOR MAKING IT NONZERO, BUT THE WD DRIVES THEM TOWARDS ZERO)

LETS CALL THE UPDATE $\Delta$. IN CASE OF SGD, $\Delta = \dfrac{\partial E}{\partial w_i}$, IN CASE OF MOMEN-TUM, $\Delta = \nu_i$.

$$\overline{w_i}^{NEW} = \overline{w_i}^{OLD} - \eta \left( \Delta + \gamma \, \overline{w_i}^{OLD} \right)$$

WEIGHT DECAY $\gamma \sim 1e^{-4}$

$\cdots 1e^{-8}$

CAN BE REORGANIZED:

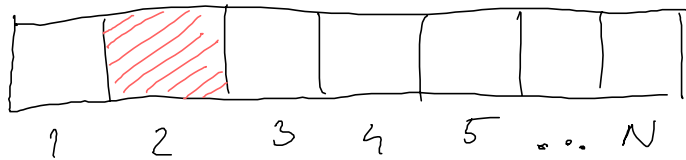$$\overline{w_i}^{NEW} = \left( 1 - \eta \gamma \right) \overline{w_i}^{OLD} - \eta \Delta$$

# CROSS VALIDATION

IN CASE THE AMOUNT OF DATA AVAILABLE IS LOW, SPLITTING THE DATASET INTO TRAIN/VALIDATION/TEST SETS MIGHT BE A PROBLEM (TOO MUCH DATA IS WASTED)

CROSS VALIDATION AVOIDS EXCLUDING THE VALIDATION DATA FROM TRAINING~ BUT THE MODEL SHOULD BE TRAINED MULTIPLE TIMES. (TAKES MUCH LONGER).

## N - WAY CROSS VALIDATION

① DIVIDE THE TRAINING SET TO N CHUNKS



1  2  3  4  5 ... N

② TRAIN N DIFFERENT MODELS:
FOR ith MODEL, EXCLUDE ith BLOCK FROM TRAINING, AND USE IT FOR VALIDATION. IF THERE ARE M SAMPLES, YOU WILL TRAIN ON $\frac{N-1}{N}M$ AND VALIDATE ON $\frac{M}{N}$ OF THEM.

3. CALCULATE YOUR VALIDATION ACCURACY BY AVERAGING THE ACCURACIES OF ALL N MODELS.

4. SELECT HYPERPARAMETERS BASED ON THE AVERAGE PERFORMANCE

5. RETRAIN THE MODEL ON FULL TRAINING DATA WITH THE SELECTED HYPER-PARAMETERS