

# MACHINE LEARNING (COURSERA)

## SUPERVISED LEARNING

$x$  := EXPLANATORY FEATURES

$y$  := TARGET FEATURE

$h$  := HYPOTHESIS FUNCTION

$J$  := COST FUNCTION

$\theta$  := PARAMETERS (BETAS)

$$\left\{ \begin{array}{l} \text{MINIMIZE } \frac{1}{2m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2 \\ \theta_0, \theta_1 \end{array} \right.$$

$$\left\{ \begin{array}{l} J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2 \\ \text{"MEAN SQUARED ERROR"} \end{array} \right.$$

$$\left\{ \begin{array}{l} h(x^{(i)}) = \theta_0 + \theta_1 x^{(i)} \end{array} \right.$$

\* LINEAR REGRESSION UNIVARIATE

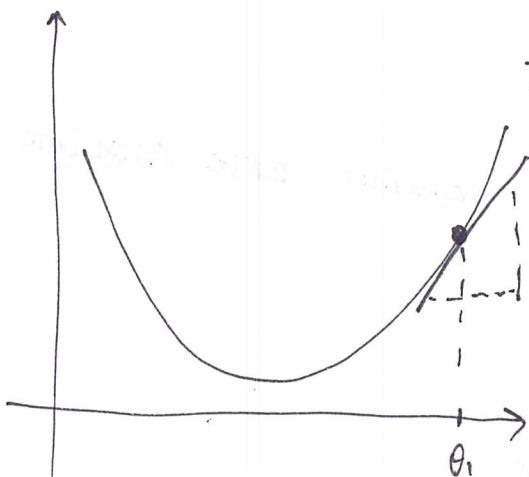
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad \text{"GRADIENT DESCENT"}$$

$$\left\{ \begin{array}{l} \text{temp}_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) \\ \text{temp}_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) \end{array} \right.$$

$$\theta_0 := \text{temp}_0$$

$$\theta_1 := \text{temp}_1$$

"BATCH": EACH STEP OF GRADIENT DESCENT USES ALL TRAINING EXAMPLES.



$$\theta_1 := \theta_1 - \underbrace{\alpha \frac{\partial}{\partial \theta_1} J(\theta_1)}_{\text{SLOPE POSITIVE } (\geq 0)}$$

$$\therefore \theta_1 := \theta_1 - \alpha \cdot (\text{POSITIVE NUMBER})$$

So  $\theta_1$  IS GOING TO THE LEFT.

\*  $\alpha$  SMALL  $\Rightarrow$  CONVERGE IS SLOW  
\*  $\alpha$  LARGE  $\Rightarrow$  PASS OVER MINIMUM OR DIVERGE!

## FOR MULTIPLE FEATURES:

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n, \quad x_0 = 1$$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}, \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \therefore \boxed{h_{\theta}(x) = \theta^T \cdot x}$$

## FEATURE SCALING

FOR FASTER PROCESS OF GRADIENT DESCENT, IT MAY BE A GOOD IDEA TO SCALE EVERY FEATURE INTO THE FOLLOWING RANGE:  $-1 \leq x_i \leq 1$

### MEAN NORMALIZATION

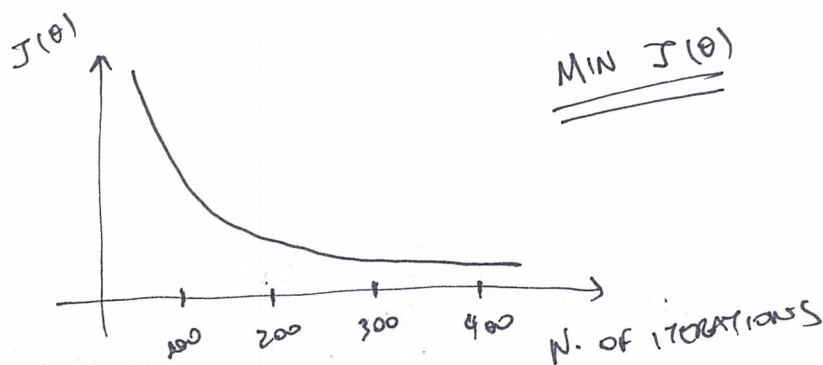
$$x_i \leftarrow \frac{x_i - \mu_1}{S_1}$$

$\mu_1 :=$  AVG VALUE OF  $x_1$  IN TRAINING SET  
 $S_1 :=$  RANGE (MAX - MIN)

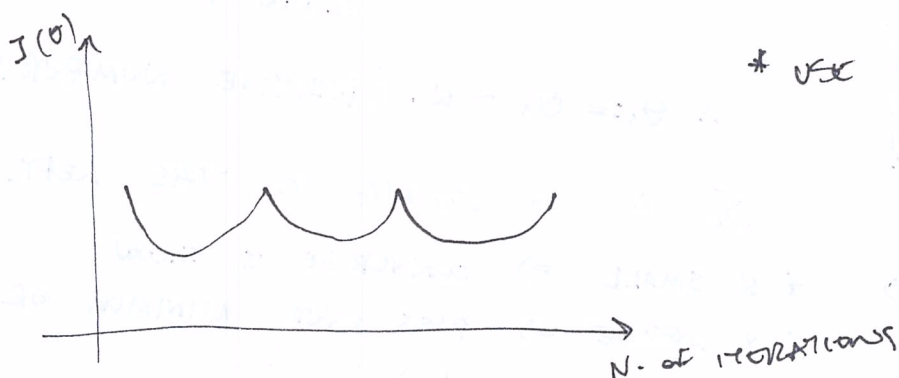
### MIN NORMALIZATION

$$x_i \leftarrow \frac{x_i - (\text{MIN } x_i)}{S_1}$$

## GRADIENT DESCENT & LEARNING RATE



\* GRADIENT DESCENT IS WORKING IF VALUE OF  $J(\theta)$  IS GETTING SMALLER AMONG THE ITERATIONS.



\* USE LEARNING RATE SMALLER ( $\alpha$ )

## NORMAL EQUATION

$n$				
INTERCEPT $x_0$	$x_1$	$x_2$	$x_3$	$y$
$m$	$\}$	$\}$	$\}$	$\}$
$\hookrightarrow X$				$\hookrightarrow y$

$$X \in \mathbb{R}^{m \times (n+1)}$$

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

$$\theta = (X^T X)^{-1} X^T y$$

## GRADIENT DESCENT

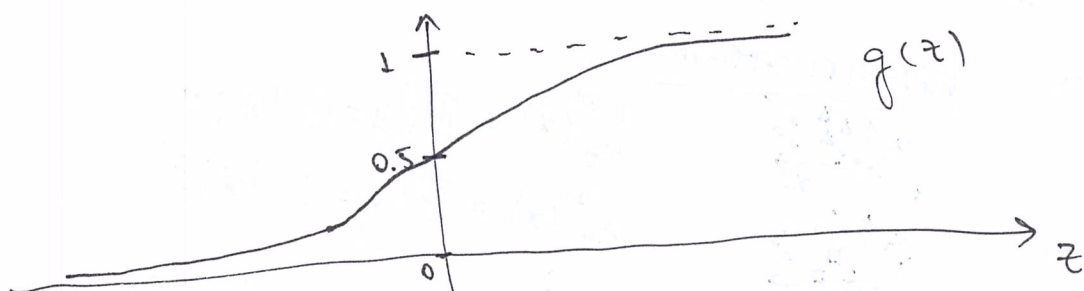
- NEED TO CHOOSE  $\alpha$
- NEEDS MANY ITERATIONS
- WORKS WELL EVEN WHEN  $n$  IS LARGE

## NORMAL EQUATION

- NO NEED TO CHOOSE  $\alpha$
- DON'T NEED TO ITERATE
- NEED TO COMPUTE  $(X^T X)^{-1}$
- SLOW IF  $n$  IS LARGE  
( $n > 10,000$ )

## LOGISTIC REGRESSION MODEL

$$0 \leq h_{\theta}(x) \leq 1, \quad h_{\theta}(x) = g(\theta^T x), \quad g(z) = \frac{1}{1 + e^{-z}} \quad (z = \theta^T x)$$



$$\therefore h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}} = \text{ESTIMATED PROBABILITY THAT } y=1 \text{ ON INPUT } x$$

SIGMOID FUNCTION  
=  
LOGISTIC FUNCTION



$h_\theta(x) = P(y=1 | x; \theta)$  "PROBABILITY THAT  $y=1$ , GIVEN  $x$ , PARAMETERIZED BY  $\theta$ "

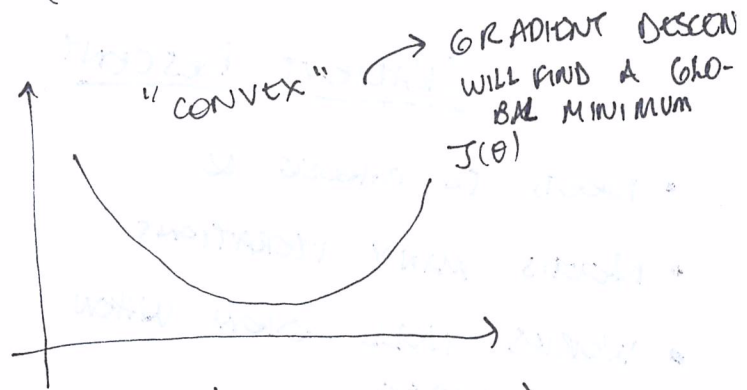
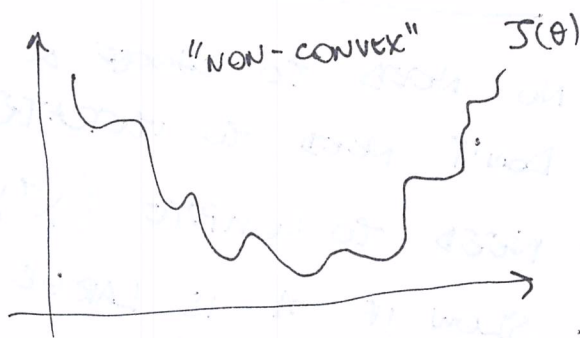
$$P(y=0 | x; \theta) + P(y=1 | x; \theta) = 1$$

## COST FUNCTION

LINEAR REGRESSION:  $J(\theta) = \frac{1}{n} \sum_{i=1}^n \underbrace{\frac{1}{2} (h_\theta(x^{(i)}) - y^{(i)})^2}_{\text{Cost}(h_\theta(x^{(i)}), y)}$

FOR LOGISTIC REGRESSION:  $h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$  NON LINEAR

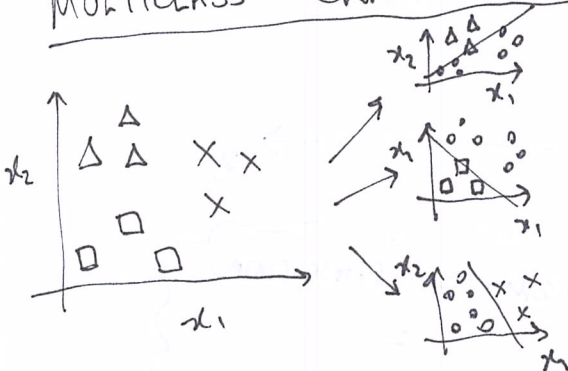
SO  $J(\theta)$  WILL BE NON-CONVEX (WITH MULTIPLE LOCAL MIN.)



LOGISTIC REGRESSION:  $\text{Cost}(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)) & , y=1 \\ -\log(1-h_\theta(x)) & , y=0 \end{cases}$

$$\therefore J(\theta) = -\frac{1}{n} \left[ \sum_{i=1}^n y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_\theta(x^{(i)})) \right]$$

## MULTICLASS CLASSIFICATION



$$h_\theta^{(0)}(x) = P(y=0 | x; \theta)$$

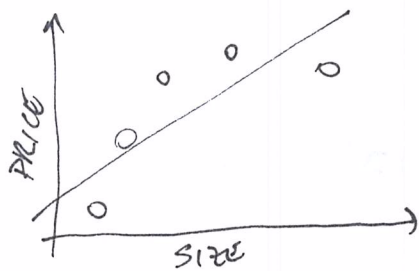
$$h_\theta^{(1)}(x) = P(y=1 | x; \theta)$$

$$\vdots$$

$$h_\theta^{(n)}(x) = P(y=n | x; \theta)$$

$$\boxed{\text{PREDICTION} = \max_i (h_\theta^{(i)}(x))}$$

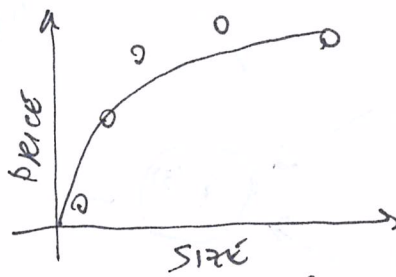
# OVERFITTING



$$\theta + \theta_1 x$$

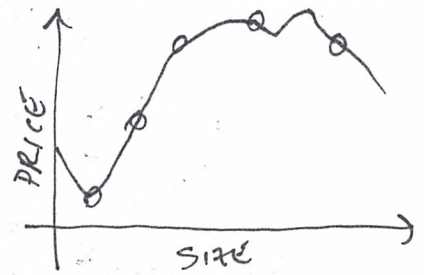
"UNDERFIT" / "HIGH BIAS"

## "LINEAR REGRESSION EXAMPLE"



$$\theta_0 + \theta_1 x + \theta_2 x^2$$

"JUST RIGHT"



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

"OVERFIT" / "HIGH VARIANCE"

DEF.: IF TOO MANY FEATURES, THE MODEL MAY FIT THE TRAINING SET TOO WELL ( $J(\theta) \approx 0$ ), BUT FAIL TO GENERALIZE TO NEW EXAMPLE

### MEANING:

- BIAS: MODEL HAS A BIAS, IN A WAY THAT DOESN'T FIT TO THE DATA (OR LOOK).
- VARIANCE: MODEL ~~TOO~~ IS VARIABLE, IN A WAY THAT IT CAN FIT PERFECTLY AT ANY DATA.

### SOLVING OVERFITTING: SOME OPTIONS

- ① REDUCE NUMBER OF FEATURES (MANUALLY / MODEL SELECTION)
- ② REGULARIZATION (REDUCES MAGNITUDE/VALUES OF PARAMETERS  $\theta$ )

### FORMULA

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2$$

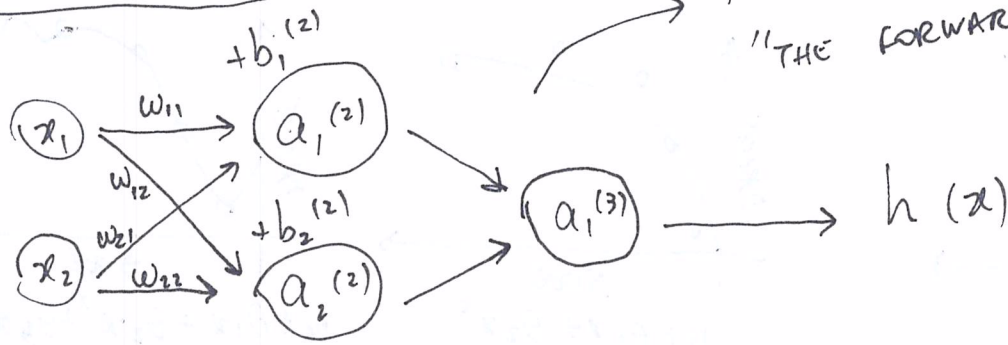
$\lambda$  := REGULARIZATION PARAMETER

REDUCES INFLUENCE OF PARAMETERS  $\theta$  (RISE COST BASED ON SUMMARIZING  $\theta$  &  $\lambda$ )

## NEURAL NETWORKS

MOTIVATION: NON-LINEAR HYPOTHESES AND NUMBER OF FEATURES ARE TOO LARGE. EX.: IMAGE RECOGNITION.  
 MOREOVER: LINEAR METHODS ARE "WEAK" (MAKE STRONG ASSUMPTIONS) AND CAN ONLY EXPRESS RELATIVELY SIMPLE FUNCTIONS OF INPUTS.

## MODEL REPRESENTATION



THIS PROCESS IS CALLED  
"THE FORWARD PROPAGATION"

LAYER 1      LAYER 2      LAYER 3  
"INPUT LAYER"    "HIDDEN LAYER"    "OUTPUT LAYER"

$$a_1^{(2)} = g(w_{11}^{(1)} \cdot x_1 + w_{21}^{(1)} \cdot x_2 + b_1^{(2)})$$

$$a_2^{(2)} = g(w_{12}^{(1)} \cdot x_1 + w_{22}^{(1)} \cdot x_2 + b_2^{(2)})$$

$$z_1^{(2)} = w_{11}^{(1)} \cdot x_1 + w_{21}^{(1)} \cdot x_2 + b_1^{(2)}$$

$a_i^{(f)} :=$  "ACTIVATION" OF UNIT  $i$  IN LAYER  $f$

$w_{ik}^{(f)} :=$  "WEIGHTS" OF LAYER  $f$  FROM UNIT  $i$  TO UNIT  $k$

$g(\cdot) :=$  ACTIVATION FUNCTION

## MULTICLASS CLASSIFICATION

$y^{(i)}$  RANGES, FOR EXAMPLE, AS:  $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$

\* KIND OF AN INTRODUCTION TO SOFTMAX

## COST FUNCTION

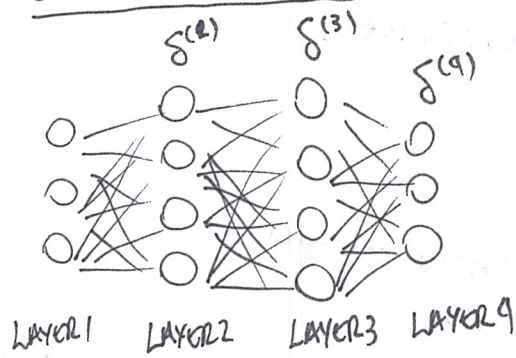
$$J(\Theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K \left[ y_k^{(i)} \log((h_\Theta(x^{(i)}))_k) + (1 - y_k^{(i)}) \log(1 - (h_\Theta(x^{(i)}))_k) \right] +$$

$$+ \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} (\Theta_{ij}^{(l)})^2$$

$$\begin{cases} L = \# \text{ TOTAL LAYERS} \\ S_l = \# \text{ UNITS IN LAYER } l \\ K = \# \text{ OUTPUT UNITS/CLASSES} \\ m = \# \text{ OF OBSERVATIONS} \end{cases}$$



## BACK PROPAGATION



$$\delta_f^{(4)} = a_f^{(4)} - y_f \quad (\text{PREDICT MINUS OBSERV.})$$

$$\delta^{(3)} = (\Theta^{(3)})^T \delta^{(4)} \cdot * g'(z^{(3)})$$

$$\delta^{(2)} = (\Theta^{(2)})^T \delta^{(3)} \cdot * g'(z^{(2)})$$

## "ALGORITHM"

TRAINING SET  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$

SET  $\Delta_{if}^{(2)} = 0$

FOR  $i = 1$  TO  $m$ :

SET  $a^{(1)} = x^{(i)}$

PERFORM FORWARD PROPAGATION TO COMPUTE  $a^{(l)}$  ( $l=2,3,\dots,L$ )

USING  $y^{(i)}$ , COMPUTE  $\delta^{(L)} = a^{(L)} - y^{(i)}$

COMPUTE  $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$

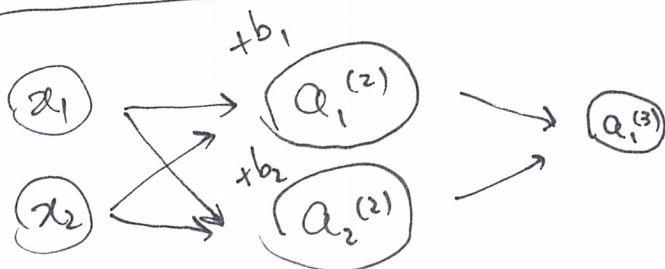
$$\Delta_{if}^{(l)} := \Delta_{if}^{(l)} + a_{if}^{(l)} \cdot \delta_i^{(l+1)}$$

$$\left\{ \begin{aligned} D_{if}^{(l)} &:= \frac{1}{m} \Delta_{if}^{(l)} + \lambda \Theta_{if}^{(l)}, & f \neq 0 \\ D_{if}^{(l)} &:= \frac{1}{m} \Delta_{if}^{(l)}, & f = 0 \end{aligned} \right.$$

(DO NOT USE REGULARIZATION ON BIAS)

$$\therefore \boxed{\frac{\partial J(\Theta)}{\partial \Theta_{if}^{(l)}} = D_{if}^{(l)}}$$

## WEIGHTS RANDOM INITIALIZATION



$$\boxed{\Theta_{if}^{(l)} = 0}, \forall i, j, l \Rightarrow$$

$$a_1^{(2)} = a_2^{(2)}, \delta_1^{(2)} = \delta_2^{(2)},$$

$$\frac{\partial}{\partial \Theta_{o1}^{(1)}} J(\Theta) = \frac{\partial}{\partial \Theta_{o1}^{(1)}} J(\Theta)$$

IF ALL HIDDEN UNITS ARE IDENTICAL THEN IT WOULD BE LIKE A MODEL OF ONLY ONE FEATURE.

## SUMMARY

- \* NUMBER OF INPUT UNITS = DIMENSION OF FEATURES  $x^{(i)}$
- \* NUMBER OF OUTPUT UNITS = NUMBER OF CLASSES
- \* NUMBER OF HIDDEN UNITS PER LAYER = USUALLY THE MORE THE BETTER
- \* NUMBER OF HIDDEN LAYERS = DEFAULT IS 1, BUT USUALLY THE MORE THE BETTER (# OF HIDDEN UNITS EQUAL BETWEEN LAYERS)

## TRAINING A NN

- ① RANDOMLY INITIALIZE THE WEIGHTS
- ② IMPLEMENT FP TO GET  $h_{\theta}(x^{(i)})$  FOR ANY  $x^{(i)}$
- ③ IMPLEMENT COST FUNCTION.
- ④ IMPLEMENT BP TO COMPUTE PARTIAL DERIVATIVES
- ⑤ USE GRADIENT CHECKING TO CONFIRM THAT BP WORKS.  
THEN DISABLE GRADIENT CHECKING.
- ⑥ USE GRADIENT DESCENT OR A BUILT-IN OPTIMIZATION FUNCTION TO MIN  $J(\theta)$  WITH WEIGHTS.