

MACHINE

LEARNING (COURSERA)

SUPERVISED LEARNING

x := EXPLANATORY FEATURES

y := TARGET FEATURE

h := HYPOTHESIS FUNCTION

J := COST FUNCTION
PARAMETERS (BETAS)

θ :=

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

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

$$\left\{ h(x^{(i)}) = \theta_0 + \theta_1 \cdot x^{(i)} \right.$$

* LINEAR REGRESSION UNIVARIATE

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

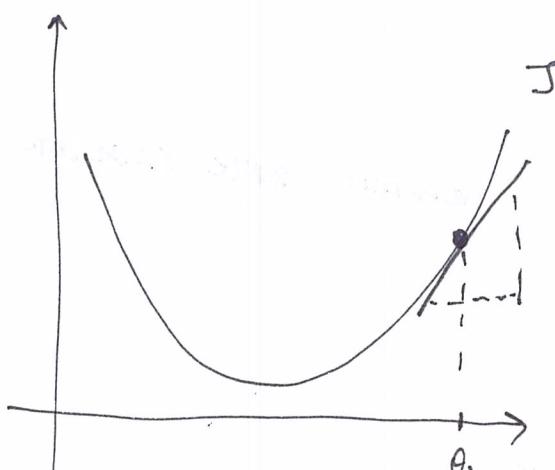
$$\text{temp_0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\text{temp_1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp_0}$$

$$\theta_1 := \text{temp_1}$$

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



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

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

So θ_1 IS GOING TO THE LEFT.

* α SMALL \Rightarrow CONVERGE IS SLOW

* α LARGE \Rightarrow PASS OVER MINIMUM OR DIVERGE!



For multiple features

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

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

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

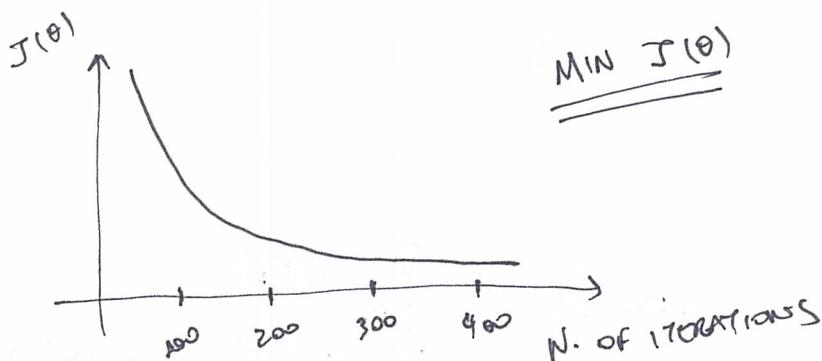
$$x_i \leftarrow \frac{x_i - \mu_i}{s_i}$$

$\mu_i :=$ avg value of x_i in training set
 $s_i :=$ range (max - min)

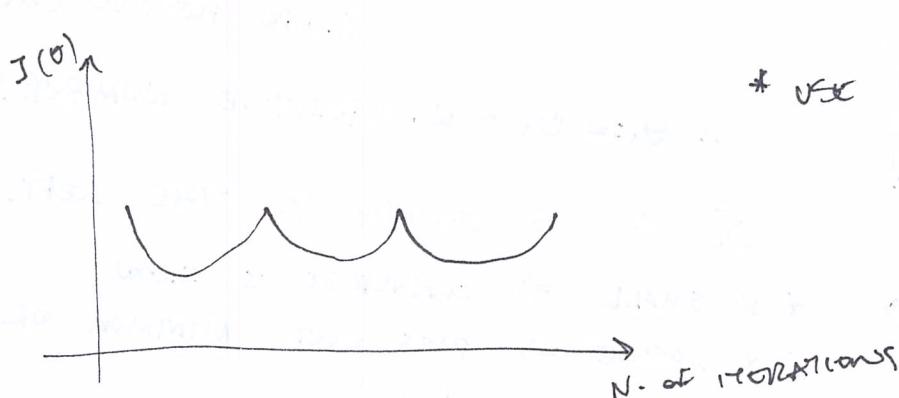
MIN NORMALIZATION

$$x_i \leftarrow \frac{x_i - (\min x_i)}{s_i}$$

GRADIENT DESCENT & LEARNING RATE



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



* USE LEARNING RATE SMALLER (α)

NORMAL EQUATION

INTERCEPT x_0	x_1	x_2	x_3	y
3	3	3	3	7

$\hookrightarrow x \quad \quad \quad \hookrightarrow y$

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

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

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

GRADIENT DESCENT

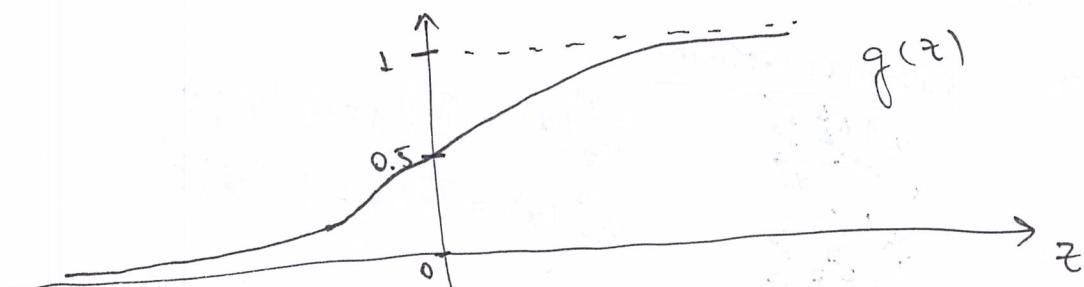
- NEEDS TO CHOOSE α
- NEEDS MANY ITERATIONS
- WORKS WELL EVEN WHEN n IS LARGE

NORMAL EQUATION

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

HOLISTIC REGRESSION MODEL

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



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

SIGMOID FUNCTION = LOGISTIC FUNCTION

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

"PROBABILITY THAT $y=1$, GIVEN x ,
PARAMETERIZED BY θ "

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

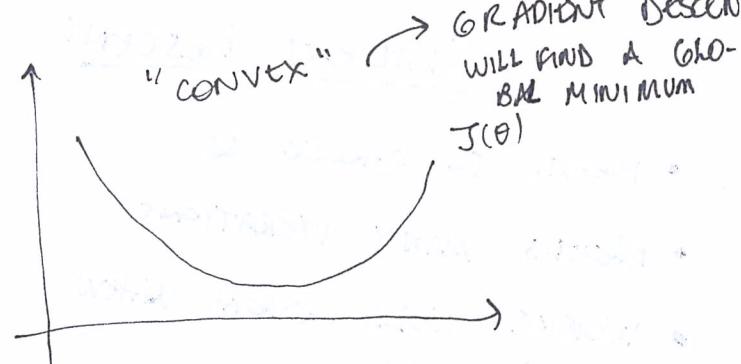
COST FUNCTION

LINAR REGRESSION : $J(\theta) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$

$\underbrace{\qquad\qquad\qquad}_{\text{cost}(h_{\theta}(x^{(i)}), y)}$

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

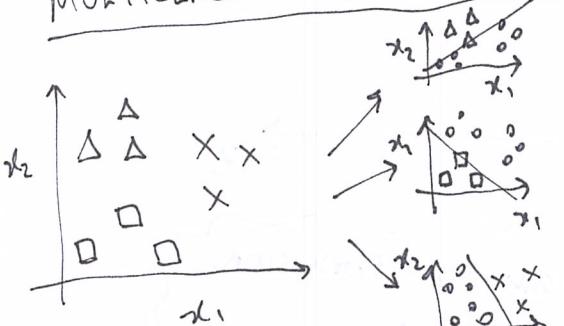
So $J(\theta)$ WILL BE NON-CONVEX (WITH MULTIPLE LOCAL MIN.)



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

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

MULTICLASS CLASSIFICATION



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

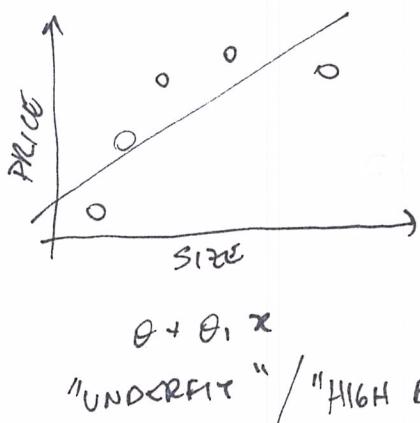
$$h_{\theta}^{(1)}(x) = P(y=1 | x; \theta)$$

$$\vdots$$

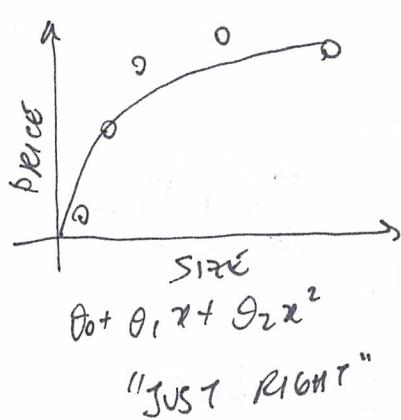
$$h_{\theta}^{(n)}(x) = P(y=n | x; \theta)$$

PREDICTION = $\max_{i=0}^n (h_{\theta}^{(i)}(x))$

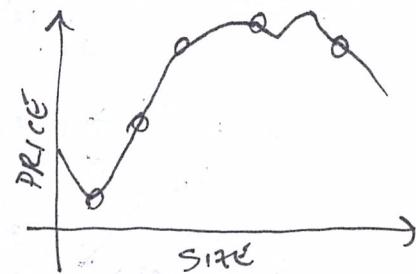
OVERRFITTING



"LINEAR REGRESSION EXAMPLE"



"JUST RIGHT"



"OVERRFIT" / "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 EXAMPLES.

MEANING: THE MODEL HAS A BIAS, IN A WAY THAT DOESN'T FIT TO THE DATA (OR LOOK).

- BIAS: MODEL HAS A BIAS, IN A WAY THAT IT CAN FIT THE DATA (OR LOOK).
- VARIANCE: MODEL IS TOO VARIABLE, IN A WAY THAT IT CAN FIT PERECTLY AT ANY DATA.

SOLVING OVERRFITTING: SOME OPTIONS

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

FORMULA

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

REDUCES INFLUENCE
OF PARAMETERS θ
(RISKS COST BASED ON
SUMMARIZING $\theta + \lambda$)

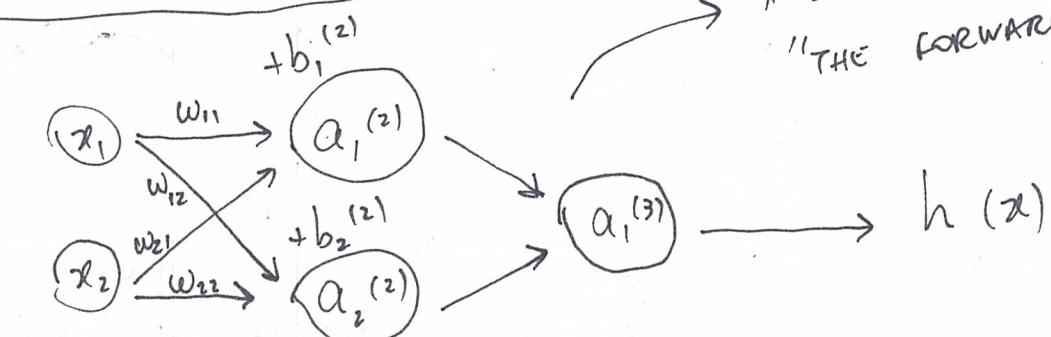
λ := REGULARIZATION PARAMETER

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



LAYER 1 LAYER 2 LAYER 3

"INPUT LAYER" "HIDDEN LAYER" "OUTPUT LAYER"

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

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

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

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

$w_{ik}^{(j)}$:= "WEIGHTS" OF LAYER j 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

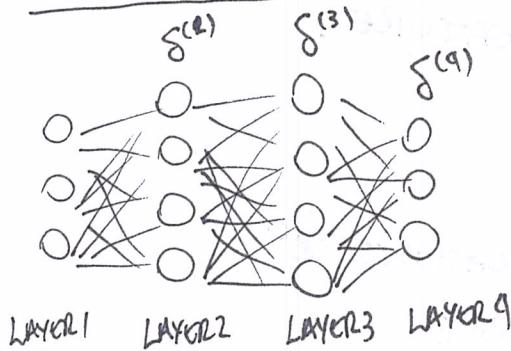
COST FUNCTION

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

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

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

BACK PROPAGATION



$$\delta_j^{(4)} = a_j^{(4)} - y_j \quad (\text{PREDICT MINUS OBSERV.})$$

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

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

"ALGORITHM"

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

$$\text{SET } \Delta_{ij}^{(l)} = 0$$

FOR $i=1$ TO m :

$$\text{SET } a^{(i)} = x^{(i)}$$

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

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

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

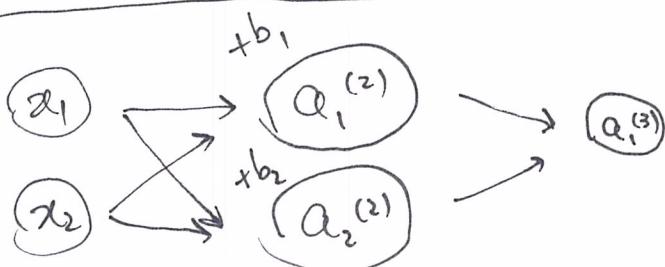
$$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_i^{(l)} \cdot \delta_j^{(l+1)}$$

$$\left\{ \begin{array}{ll} D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)} & , j \neq 0 \\ D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} & , j = 0 \end{array} \right.$$

(DO NOT USE REGULARIZATION ON BIAS)

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

WEIGHTS RANDOM INITIALIZATION



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

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

$$\frac{\partial}{\partial \Theta_{01}^{(1)}} J(\Theta) = \frac{\partial}{\partial \Theta_{02}^{(1)}} J(\Theta)$$

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

SUMMARY

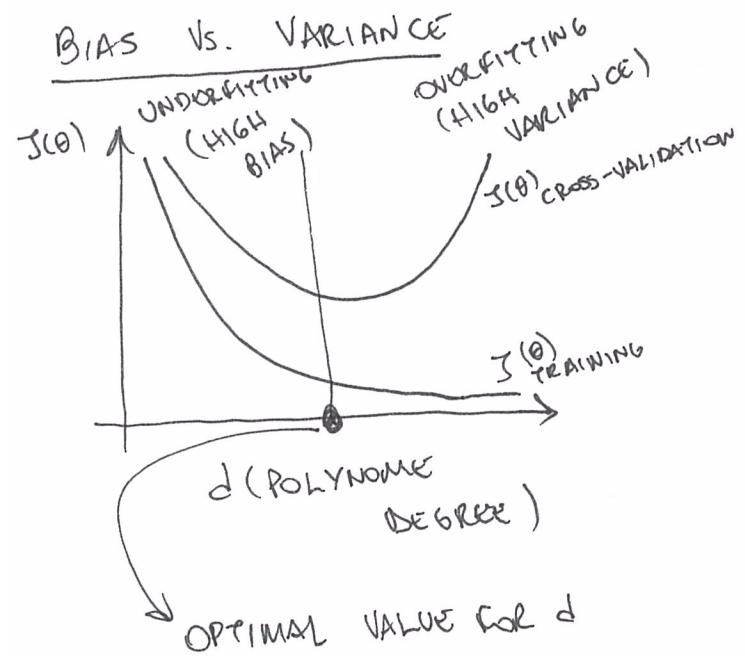
- * NUMBER OF INPUT UNITS = DIMENSION OF FEATURES $x^{(i)}$
- * NUMBER OF OUTPUT UNITS = NUMBER OF CLASSES
- * NUMBER OF HIDDEN UNITS PER LAYER = USUALLY THE MORE THE BETTER
- * NUMBER OF HIDDEN LAYERS = DEFAULT IS 1, BUT USUALLY THE MORE THE BETTER (NUMBER OF HIDDEN UNITS EQUAL BETWEEN LAYERS)

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.
- ⑥ THEN DISABLE GRADIENT CHECKING.
USE GRADIENT DESCENT OR A BUILT-IN OPTIMIZATION FUNCTION TO MIN $J(\theta)$ WITH WEIGHTS.

MODEL SELECTION

- ① OPTIMIZE THE PARAMETERS θ (ESTIMATIONS) USING THE TRAINING SET FOR EACH POLYNOMIAL DEGREE
- ② FIND THE POLYNOMIAL DEGREE d WITH THE LEAST ERROR USING CROSS-VALIDATION SET
- ③ ESTIMATE FINAL ERROR USING TEST SET.



REGULARIZATION

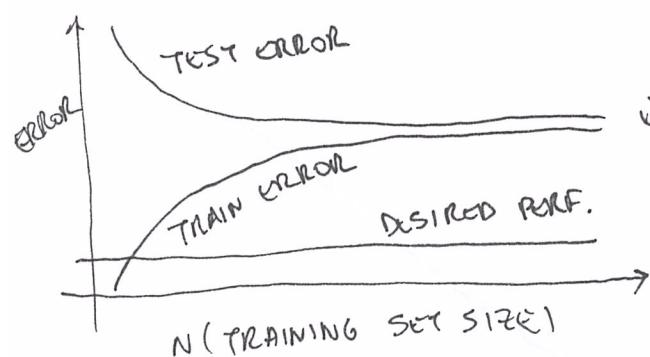
$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{m} \sum_{j=1}^n \theta_j^2$$

- LARGE $\lambda \Rightarrow$ HIGH BIAS (UNDERRFIT)
 $\lambda = 10000 \Rightarrow \theta_1 \approx 0, \theta_2 \approx 0 \dots h(x) \approx \theta_0$
- INTERMEDIATE $\lambda \Rightarrow$ JUST RIGHT
- SMALL $\lambda \Rightarrow$ HIGH VARIANCE (OVERFIT)
 $\lambda = 0$

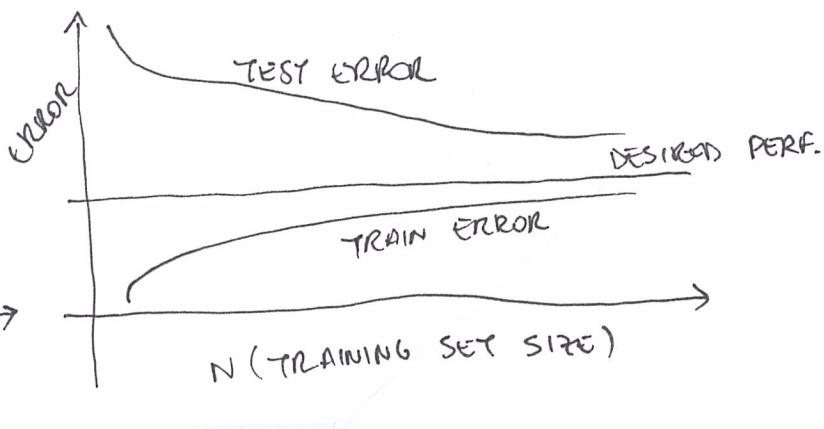
ITERATE THROUGH LIST OF LAMBDRAS:
 $\lambda \in \{0, 0.01, 0.02, 0.08, 0.16, 0.32, 0.64, 1.28, 2.56, 5.12, 10.24\}$

LEARNING CURVES

HIGH BIAS:



HIGH VARIANCE:



SUMMARY

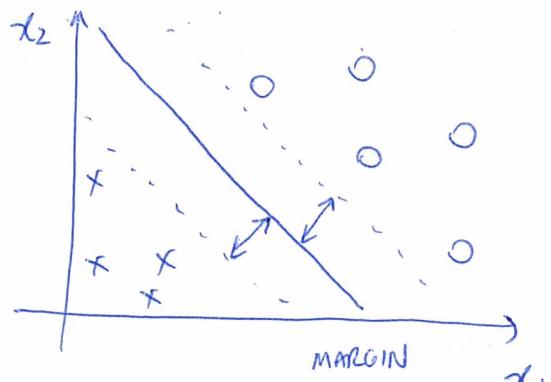
- GETTING MORE TRAINING EXAMPLES: FIXES HIGH VARIANCE
- TRYING SMALLER SETS OF FEATURES: FIXES HIGH VARIANCE
- ADDING FEATURES: FIXES HIGH BIAS
- ADDING POLYNOMIAL FEATURES: FIXES HIGH BIAS
- DECREASING λ : FIXES HIGH BIAS
- INCREASING λ : FIXES HIGH VARIANCE

SUPPORT VECTOR MACHINE

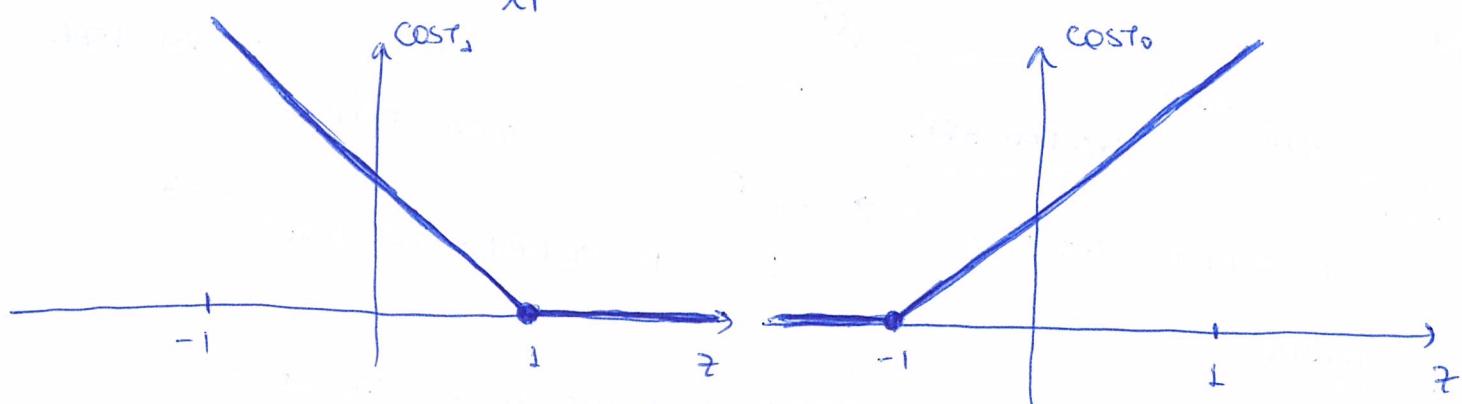
$$J(\theta) = \min_{\theta} C \cdot \sum_{i=1}^m \left[y^{(i)} \text{cost}_1 (\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0 (\theta^T x^{(i)}) \right] +$$

$$\frac{1}{2} \sum_{j=1}^n \theta_j^2 \quad (\text{COST FUNCTION OF SVM})$$

s.t. $\begin{cases} \theta^T x^{(i)} \geq 1, & \text{if } y^{(i)} = 1 \\ \theta^T x^{(i)} \leq -1, & \text{if } y^{(i)} = 0 \end{cases}$

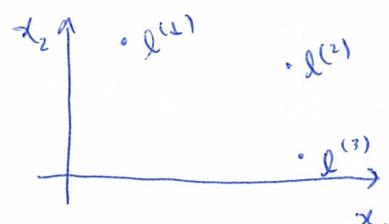
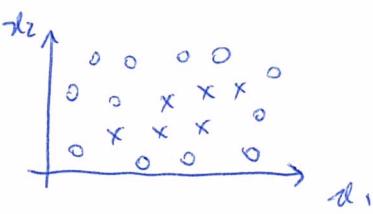


LINER DECISION BOUNDARY



KERNELS

(FOR NON-LINEAR DECISION BOUNDARY)



$$f_i = \text{SIMILARITY } (x, l^{(i)}) \\ = \exp \left(- \frac{\|x - l^{(i)}\|^2}{2 \sigma^2} \right)$$

"GAUSSIAN KERNEL"

$$J(\theta) = \min_{\theta} C \cdot \sum_{i=1}^m \left[y^{(i)} \text{cost}_1 (\theta^T f_i) + (1 - y^{(i)}) \text{cost}_0 (\theta^T f_i) \right] +$$

$$\frac{1}{2} \sum_{j=1}^n \theta_j^2 \quad (\text{COST FUNCTION OF SVM FOR KERNELS})$$

SVM PARAMETERS

$C \left(= \frac{1}{\lambda}\right)$ } LARGE : LOW BIAS, HIGH VARIANCE
SMALL : HIGH BIAS, LOW VARIANCE

σ^2 } LARGE : FEATURES f_i VARY MORE SMOOTHLY. HIGH BIAS, LOW VARIANCE
SMALL : FEATURES f_i VARY LESS SMOOTHLY. LOW BIAS, HIGH VARIANCE

MULTI-CLASS CLASSIFICATION

TRAIN \leq SVMs, ONE TO DISTINGUISH $y = i$ FROM THE REST
for $i = 1, 2, \dots, K$: $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(K)}$
PICK CLASS i WITH LARGEST $(\theta^{(i)})^T x$

LOGISTIC REG. VS. SVMs

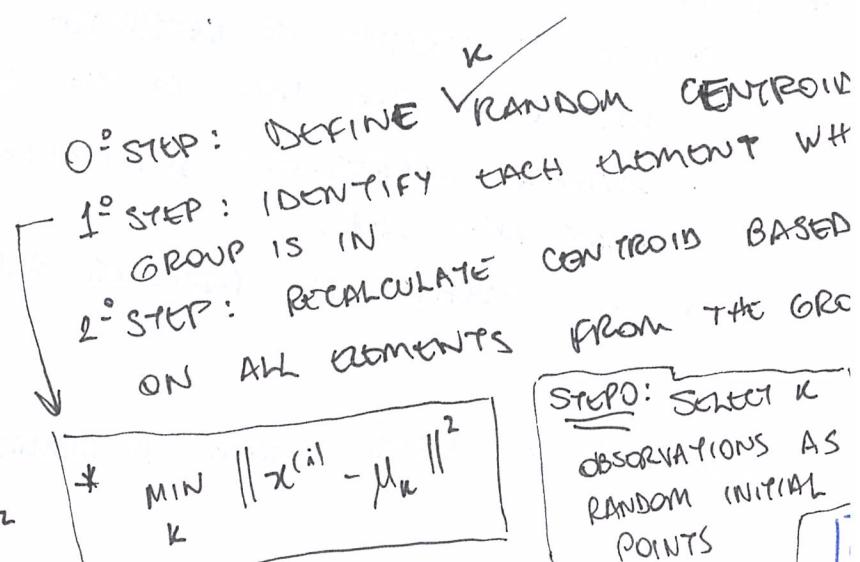
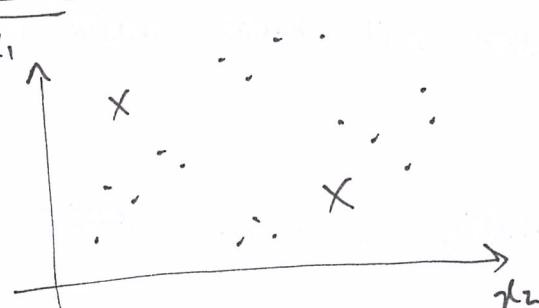
n = NUMBER OF FEATURES
 m = NUMBER OF TRAINING EXAMPLES

- n LARGE : LOGISTIC REGRESSION OR SVM WITHOUT KERNEL
- n SMALL, m INTERMEDIATE : SVM WITH GAUSSIAN KERNEL
- n SMALL, m LARGE : LOGISTIC REG. OR SVM WITHOUT A KERNEL

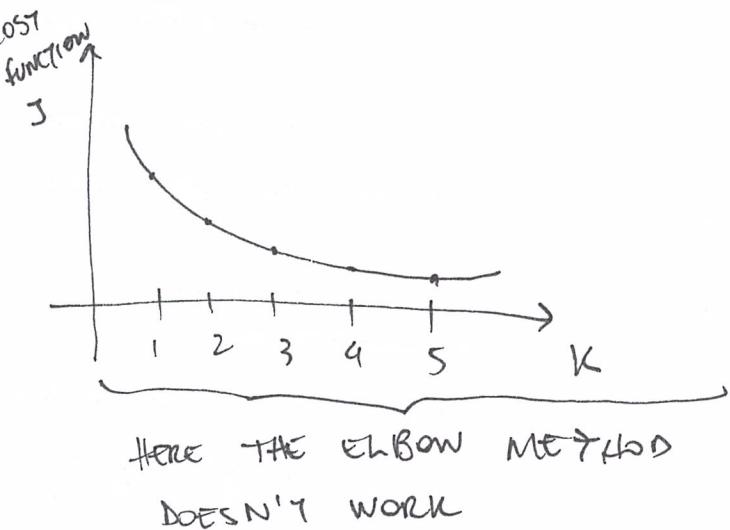
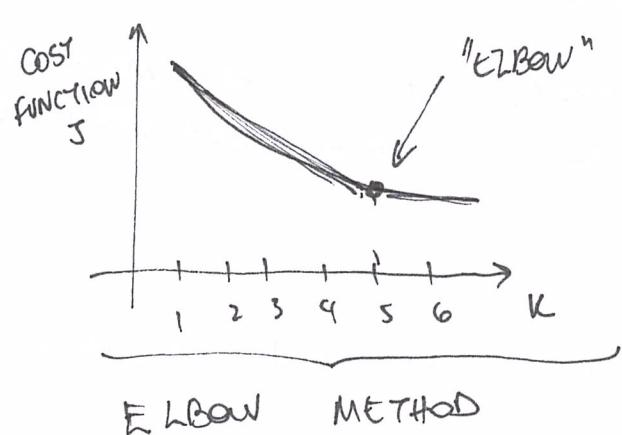
UNSUPERVISED LEARNING

GIVEN AN UNLABELED DATASET, FIND "STRUCTURE" IN THE DATA. TRAINING SET OF FORM $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ WITHOUT LABELS $y^{(i)}$.

K-MEANS

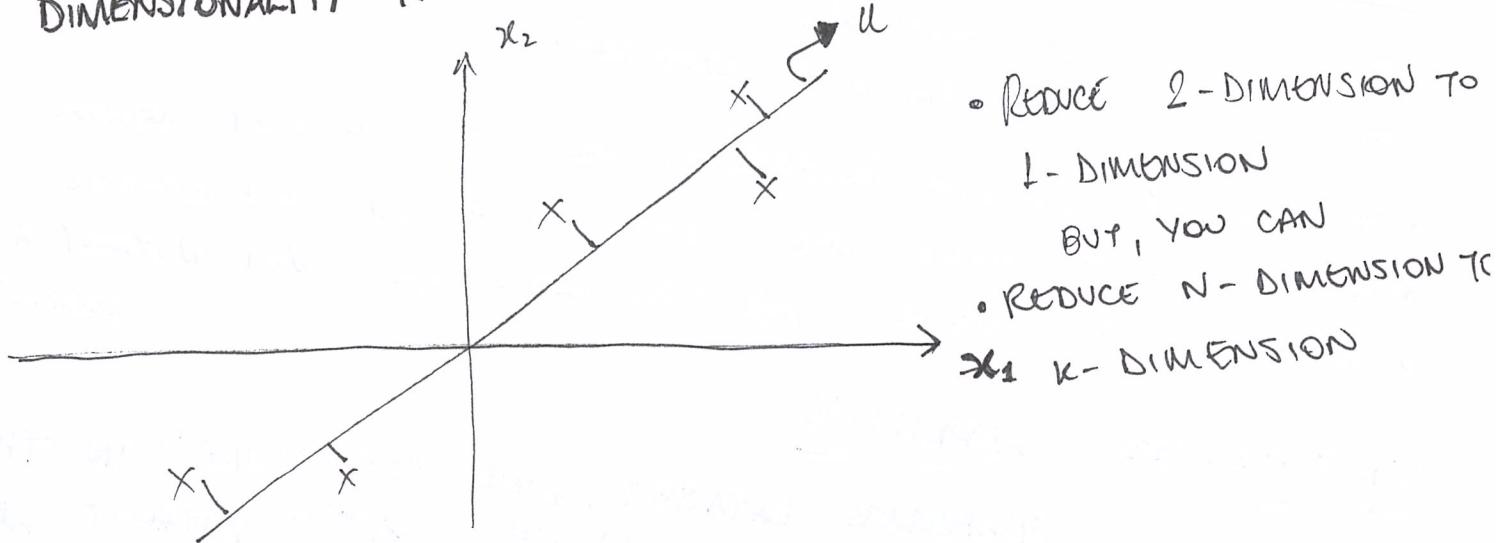


CHOOSING THE VALUE OF K



PCA : PRINCIPAL COMPONENT ANALYSIS

DIMENSIONALITY REDUCTION : DATA COMPRESSION & VISUALIZATION



- * PCA IS NOT LINEAR REGRESSION
- PCA IS TRYING TO OPTIMIZE THE SHORTEST DISTANCE ORTHOGONAL FROM THE DOTS TO THE LINE.
- LINEAR REGRESSIONS MINIMIZE THE SQUARE ERROR BETWEEN POINT AND LINE.
- PCA TREATS ALL FEATURES ALIKE, DIFF. FROM LINEAR REG.

To REDUCE DATA FROM n -DIMENSIONS TO k -DIMENSIONS :

$$\bullet \quad \sum = \frac{1}{m} \sum_{i=1}^n \underbrace{(x^{(i)})}_{n \times 1} \underbrace{(x^{(i)})^T}_{1 \times n}$$

"COVARIANCE MATRIX" = SIGMA

• SVD (SIGMA) $\rightarrow U$ "EINVECTORS" of MATRIX \sum = SINGULAR VALUE DECOMPOSITION

$$U = \begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ u^{(1)} & u^{(2)} & u^{(3)} & \dots & u^{(n)} \\ 1 & 1 & 1 & \dots & 1 \end{bmatrix}, \quad U \in \mathbb{R}^{n \times n}$$

$\underbrace{\qquad\qquad\qquad}_{K}$

$$Z = \begin{bmatrix} 1 & 1 & \dots & 1 \\ u^{(1)} & u^{(2)} & \dots & u^{(n)} \\ 1 & 1 & \dots & 1 \end{bmatrix}^T X = \begin{bmatrix} -(u^{(1)})^T \\ -(u^{(2)})^T \\ \vdots \\ -(u^{(n)})^T \end{bmatrix} X$$

$\underbrace{\qquad\qquad\qquad}_{n \times K}$ $\underbrace{\qquad\qquad\qquad}_{K \times n}$ $\underbrace{\qquad\qquad\qquad}_{K \times 1}$

PCA TRIES TO MINIMIZE: $\frac{1}{m} \sum_{i=1}^n \|x^{(i)} - \hat{x}^{(i)}_{\text{APPROX}}\|^2$

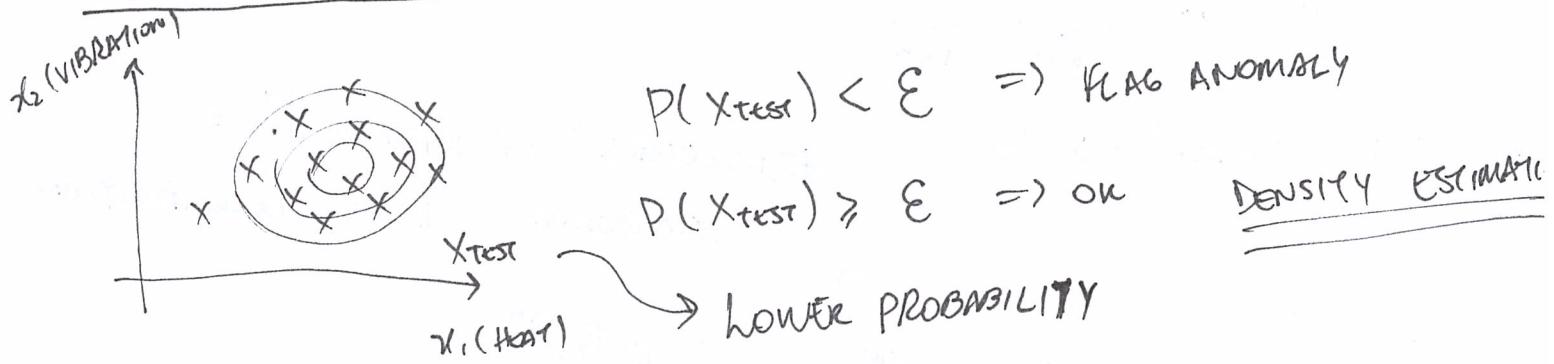
WHERE $\hat{x}^{(i)}_{\text{APPROX}}$ IS THE RESULT OF THE APPLICATION OF PCA.

$$\left(\underbrace{\hat{x}_{\text{APPROX}}}_{n \times 1} = \underbrace{U_{\text{reduce}}}_{n \times K} \cdot \underbrace{Z}_{K \times 1} \right)$$

SUMMARY (WHY TO USE PCA?)

- COMPRESS THE DATA TO TAKE LESS COMPUTER MEMORY
- REDUCE DIMENSION OF INPUT DATA TO SPEED UP LEARNING
- ALGORITHM
- TO VISUALIZE HIGH-DIMENSIONAL DATA ($K=2$ OR $K=3$)

ANOMALY DETECTION



ANOMALY DETECTION vs. SUPERVISED LEARNING

- VERY SMALL N° OF POSITIVE EXAMPLES
- LARGE N° OF NEGATIVE EXAMPLES
- MANY DIFF. TYPES OF ANOMALIES
- LARGE N° OF POSITIVE AND NEGATIVE EXAMPLES
- ENOUGH POSITIVE EXAMPLES FOR ALGORITHM GET SENSE OF WHAT POSITIVE EXAMPLES ARE LIKE.

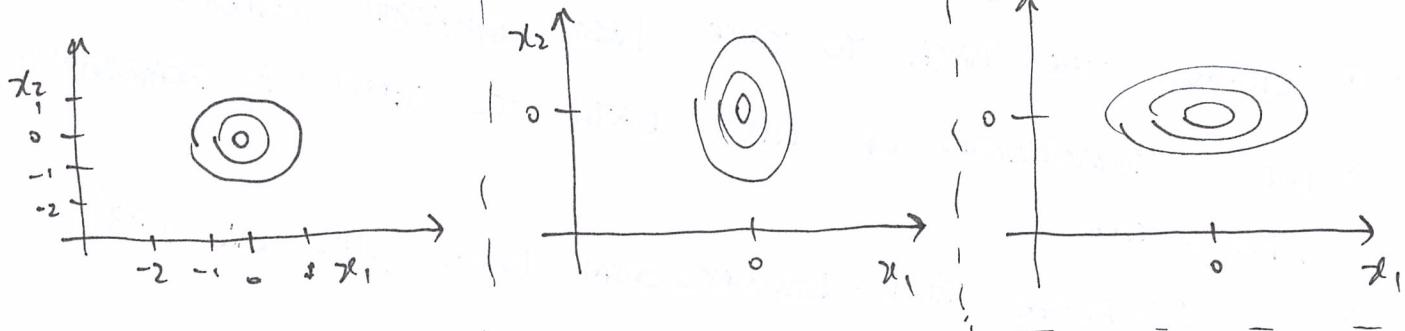
MULTIVARIATE GAUSSIAN (NORMAL) DISTRIBUTION

Don't model $p(x_1), p(x_2), \dots$ SEPARATELY, MODEL $p(x)$ ALL IN ONE:

$\mu \in \mathbb{R}^n, \Sigma \in \mathbb{R}^{n \times n}$ (COVARIANCE MATRIX)

$$p(x|\mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad | \quad \mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 1 \end{bmatrix} \quad | \quad \mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$$



CONTENT BASED RECOMMENDATIONS

MOVIE	U1	U2	U3	U4	χ_1 (ROMANCE)	χ_2 (ACTION)
A	5	5	0	0	0.9	0
B	5	?	?	0	1.0	0.01
C	?	4	0	4	0.99	0
D	0	0	5	?	0.1	-1.0

$$\chi^{(1)} = \begin{bmatrix} 1 \\ 0.9 \\ 0 \end{bmatrix}, \quad \chi^{(2)} = \begin{bmatrix} 1 \\ 1.0 \\ 0.01 \end{bmatrix}, \quad \dots$$

FEATURES FOR EACH MOVIE BASED ON USER EXPERIENCES / RATINGS

$$\theta^{(1)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}, \quad \dots$$

PARAMETER OF PREFERENCE FOR EACH USER (ESTIMATED)

PROBLEM FORMULATION

$$\min_{\theta^{(1)}} \frac{1}{2m^{(1)}} \sum_{i: r(i, 1) = 1} \left((\theta^{(1)})^T \chi^{(i)} - y^{(i, 1)} \right)^2 + \frac{\lambda}{2m^{(1)}} \sum_{k=1}^n (\theta_k^{(1)})^2$$

- $r(i, 1) = 1$, IF USER 1 HAS RATED MOVIE i (0 OTHERWISE)

- $y^{(i, 1)}$, IF RATING BY USER 1 ON MOVIE i NOT DEFINED

- $\theta^{(1)}$ PARAMETER VECTOR FOR USER 1

- $\chi^{(i)}$ FEATURE VECTOR FOR MOVIE i

- PREDICTED RATING: $(\theta^{(1)})^T \chi^{(i)}$ FOR USER 1 AND MOVIE i

COLLABORATIVE FILTERING (PROBLEM FORMULATION TO LEARN $\chi^{(i)}$)

$$\min_{\chi^{(i)}} \frac{1}{2} \sum_{j: r(i, j) = 1} \left((\theta^{(1)})^T \chi^{(i)} - y^{(i, j)} \right)^2 + \frac{\lambda}{2} \sum_{k=1}^n (\chi_k^{(i)})^2$$

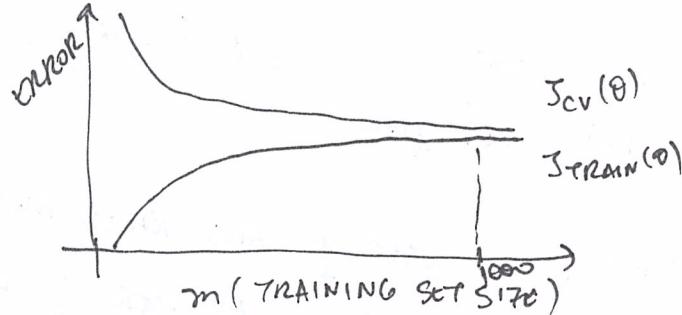
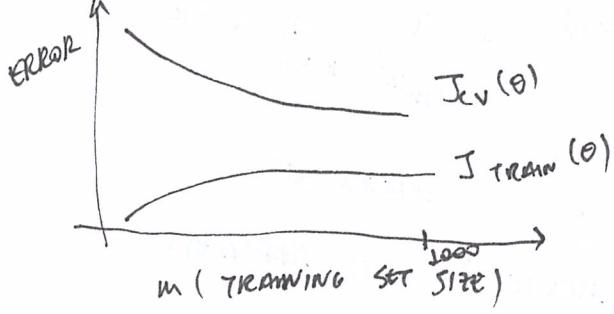
COLLABORATIVE FILTERING ALGORITHM

- ① INITIALIZE $x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}$ TO SMALL RANDOM VALUES
- ② MINIMIZE $J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)})$ USING GRADIENT DESCENT - E.G.:

$$\left\{ \begin{array}{l} x_k^{(i)} = x_k^{(i)} - \alpha \left(\sum_{j: r(x_{ij})=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) \theta_k^{(j)} + \lambda x_k^{(i)} \right) \\ \theta_k^{(j)} = \theta_k^{(j)} - \alpha \left(\sum_{i: r(x_{ij})=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \theta_k^{(j)} \right) \end{array} \right.$$

LEARNING WITH LARGE DATASETS

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$



ON THESE CASES, THE FIRST INDICATES THAT m SHOULD BE HIGHER THAN 1000, AND THE SECOND IT SEEMS THAT 1000 IS SUFFICIENT. (ERROR-DRIVEN)

- BATCH GRADIENT DESCENT: USE ALL m EXAMPLES IN EACH ITERATION
- STOCHASTIC GRADIENT DESCENT: USE 1 EXAMPLE IN EACH ITERATION
- MINI-BATCH GRADIENT DESCENT: USE \underline{B} EXAMPLES IN EACH ITERATION

ONLINE LEARNING

Given some characteristics x want to learn to predict y based on parameters $\theta := p(y=1|x; \theta)$

REPEAT forever {

GET (x, y) corresponding to user

UPDATE θ using (x, y) :

$$\theta_j := \theta_j - \alpha (h_\theta(x) - y) \cdot x_j \quad (j = 0, \dots, n)$$

}

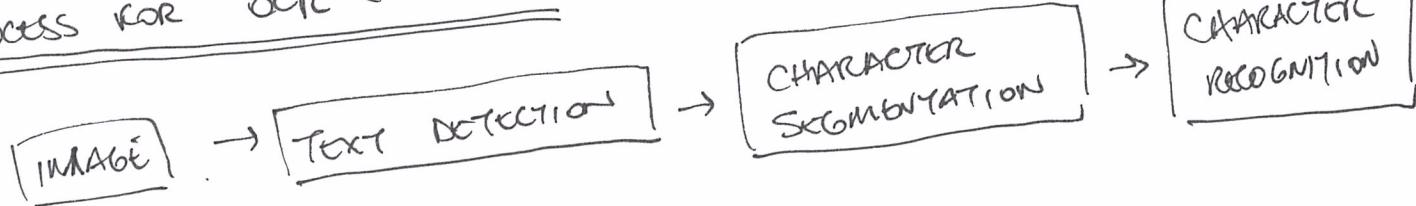
for each observation RETRAIN THE PARAMETERS θ .

MACHINE LEARNING PIPELINE & PHOTO OCR

System with many stages / components, several of which may use a ML technique.

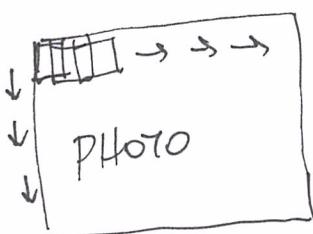
ML TECHNIQUE:
OCR: OPTICAL CHARACTER RECOGNITION

PROCESS FOR OCR (EXAMPLE):



SLIDING WINDOWS

TO COMPLETE THE FIRST STEP (TEXT DETECTION) IT'S USED THE TECHNIQUE OF SLIDING WINDOWS:



Each rectangle tries to identify if there are any type of text on it (SUPERVISED LEARNING) and after each iteration changes size of rectangle.