

ML WEEK 3: LOGISTIC REGRESSION

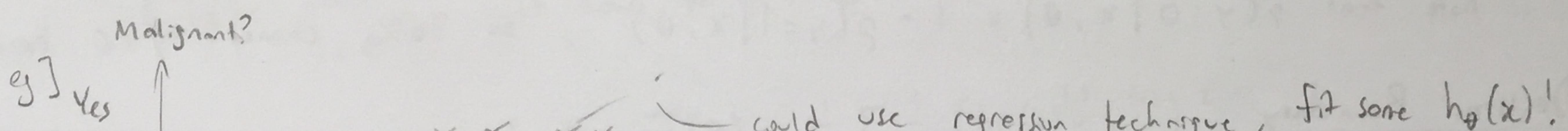
LEO ROBINOVITCH

- Classification Problems \Rightarrow spam/not spam? Fraud/Not Fraud?

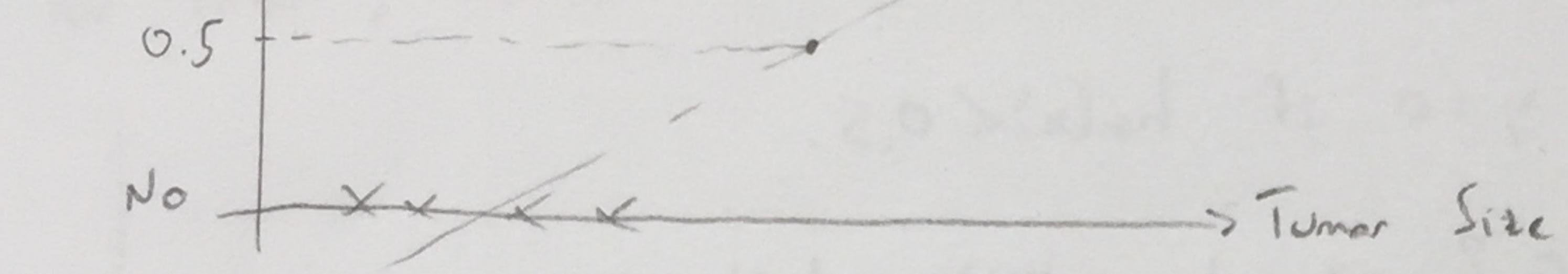
$\hookrightarrow y \in \{0, 1\} \rightarrow$

- 0: negative class (e.g. benign tumor)
- 1: positive class (e.g. malignant tumor)

MORE CLASSES POSSIBLE



could use regression technique, fit some $h_\theta(x)$!



$$\left. \begin{array}{l} \text{W/ regression, say } h_\theta(x) \geq 0.5 \rightarrow y=1 \\ h_\theta(x) < 0.5 \rightarrow y=0 \end{array} \right\}$$

OFTEN BAD

IDEA!

$\Rightarrow h_\theta(x)$ not great predictor all the time.

$\hookrightarrow h_\theta(x) > 1$ & $h_\theta(x) < 0$, too...

Logistic Regression: Hypothesis Representation:

\Rightarrow previously, $h_\theta(x) = \theta^T x$ where $x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} = 1$

\hookrightarrow in linear regression.

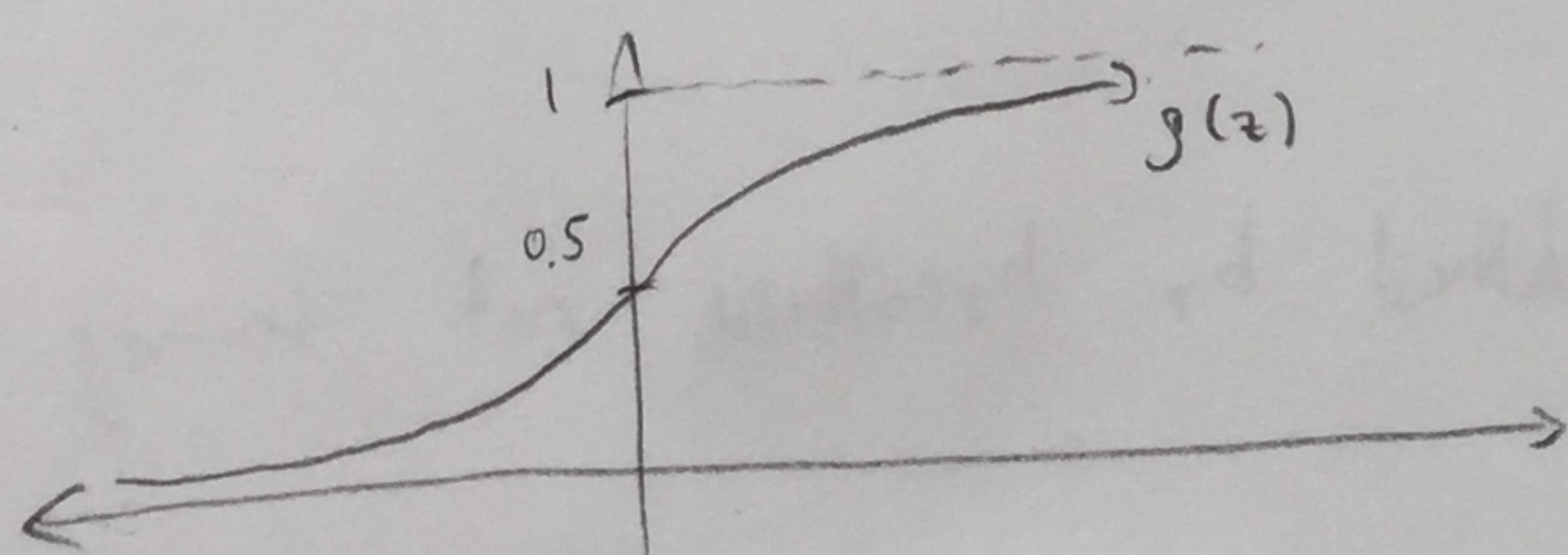
$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$, both $n+1 \times 1$ matrices where n is # features

Logistic Regression \rightarrow Sigmoid Function:

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

$$\hookrightarrow h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

$g(z)$ sigmoid or logistic fnt.



Interpretation of $h_\theta(x) \rightarrow$ estimated probability that $y=1$ on input x .

↳ e.g. $x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumor size} \end{bmatrix} \rightarrow h_\theta(x) = 0.7$, 70% chance tumor malignant

↳ $h_\theta(x) = p(y=1|x; \theta)$ "probability that $y=1$, given x , parameterized by θ "

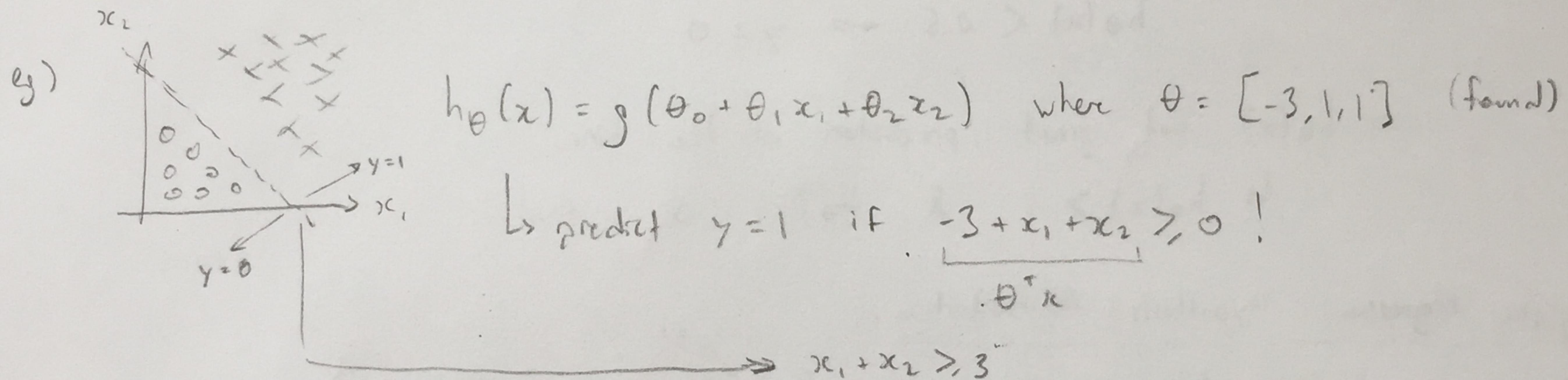
* note that $p(y=0|x; \theta) = 1 - p(y=1|x; \theta)$, so 30% chance $y=0$ here

Decision Boundary:

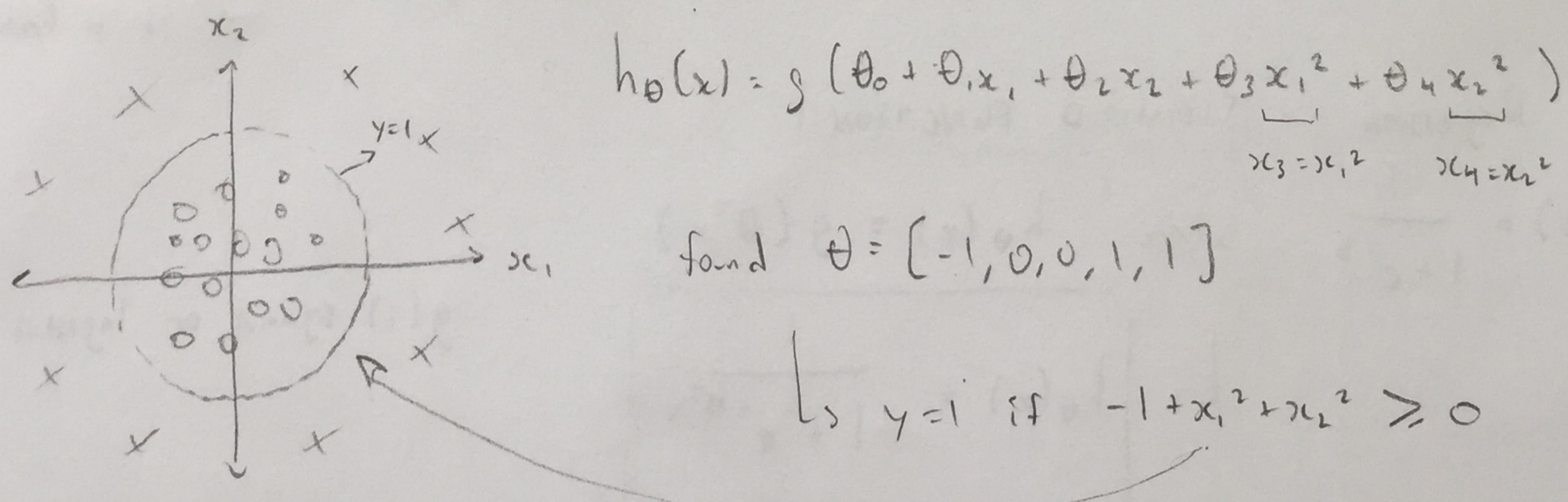
Suppose $y=1$ if $h_\theta(x) \geq 0.5$, $y=0$ if $h_\theta(x) < 0.5$.

↳ when will $h_\theta(x)$ be ≥ 0.5 ? \rightarrow when $z \geq 0$!

$$z \geq 0 \rightarrow g(\theta^T x) \geq 0.5 \text{ when } \underline{\theta^T x \geq 0}$$



Non-Linear Decision Boundaries:



→ decision boundary is defined by hypothesis, not training set!

\rightarrow Function \rightarrow getting θ

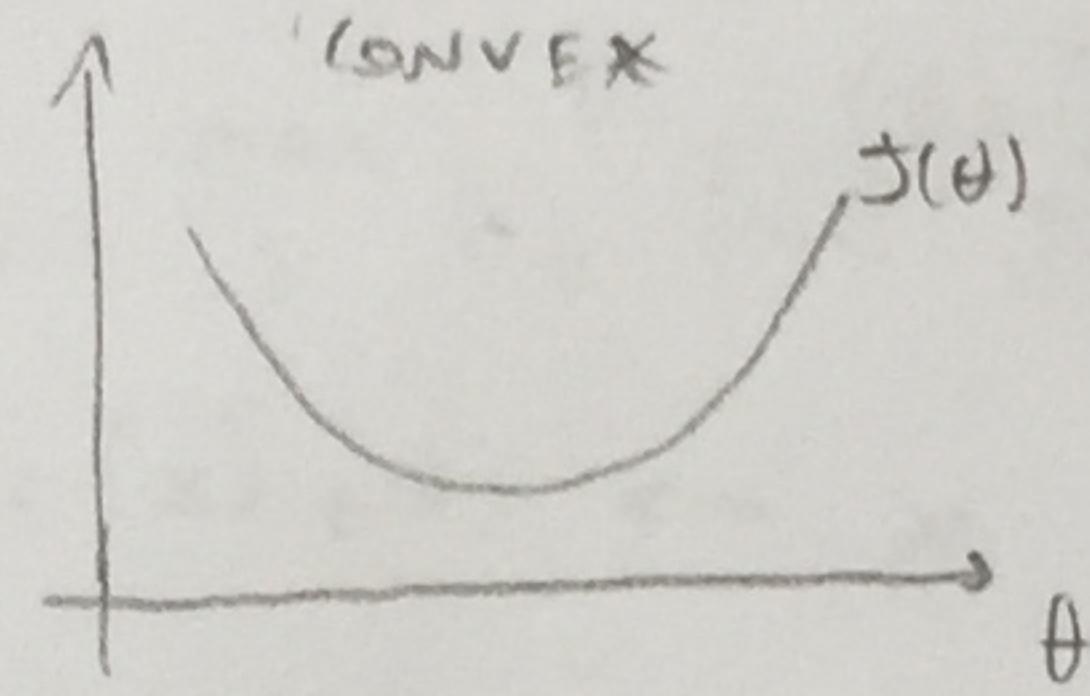
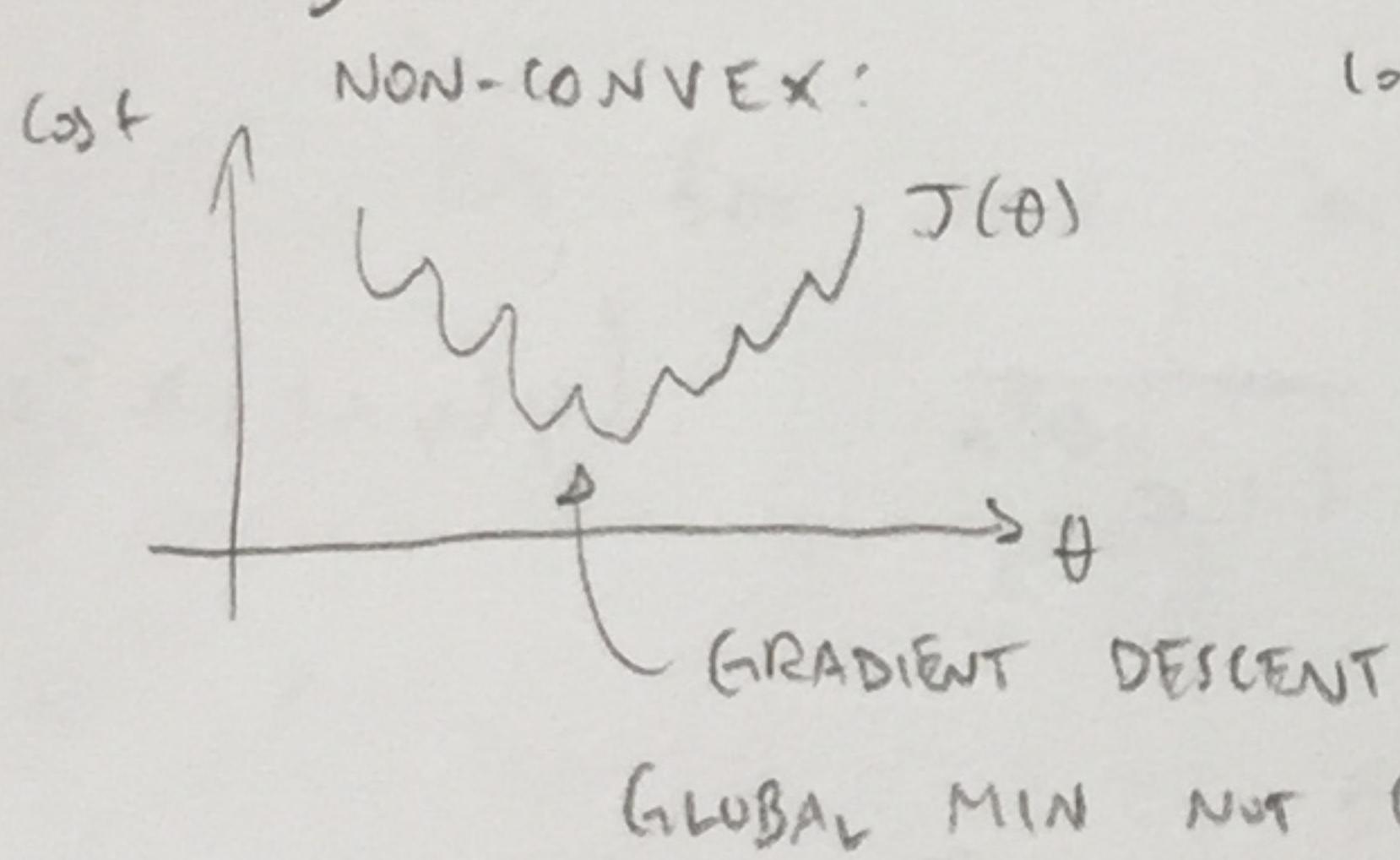
ex) Training set (m -examples) $\rightarrow \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}) \dots (x^{(m)}, y^{(m)})\}$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \text{ where } n \text{ features, } y \in \{0, 1\}$$

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}} \rightarrow \text{how to get } \theta?$$

$$\rightarrow \text{for linear, used cost function } J(\theta) = \frac{1}{m} \sum \underbrace{\frac{1}{2} (h_{\theta}(x) - y)^2}_{\text{cost}(h_{\theta}(x) - y)}$$

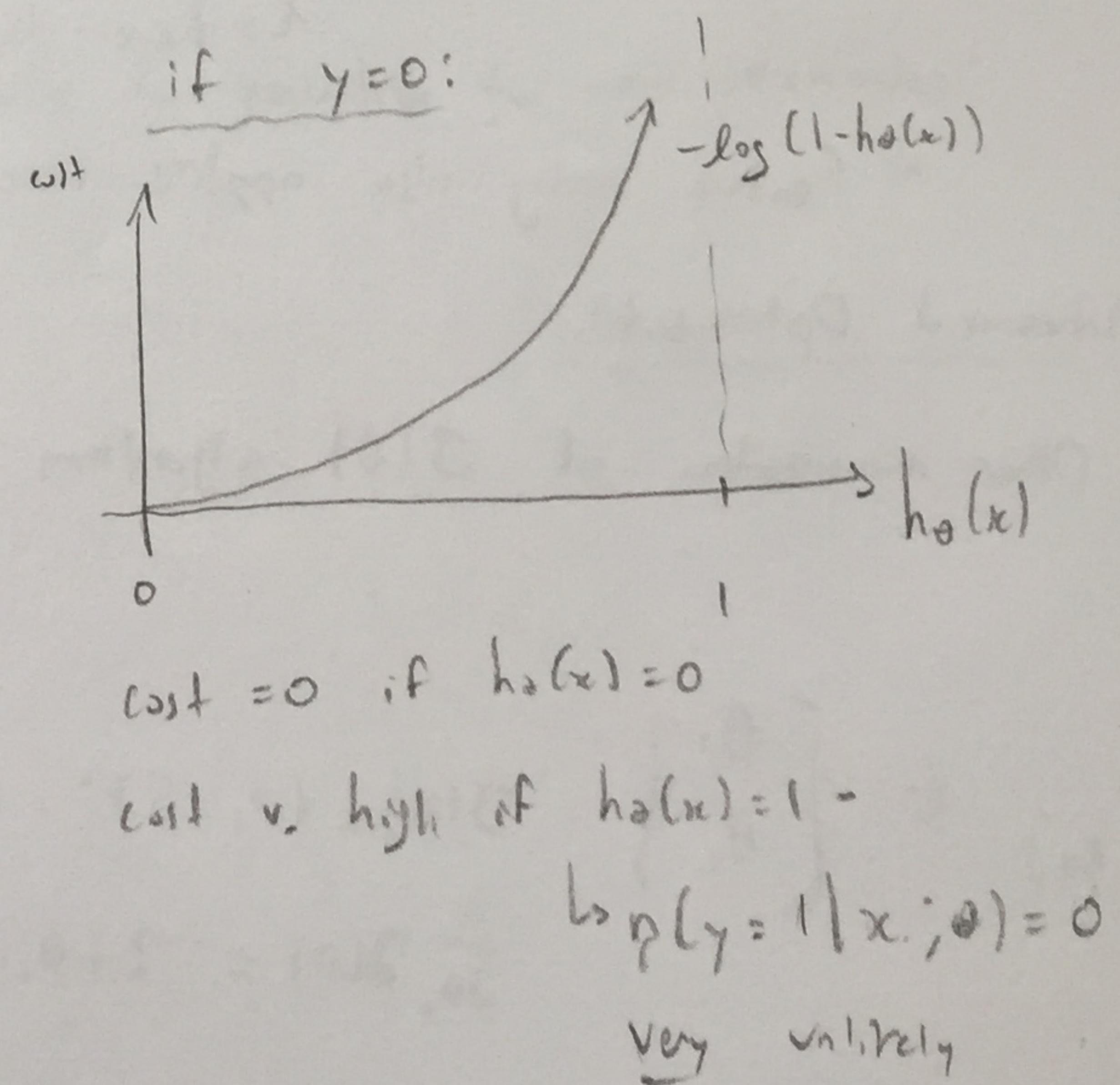
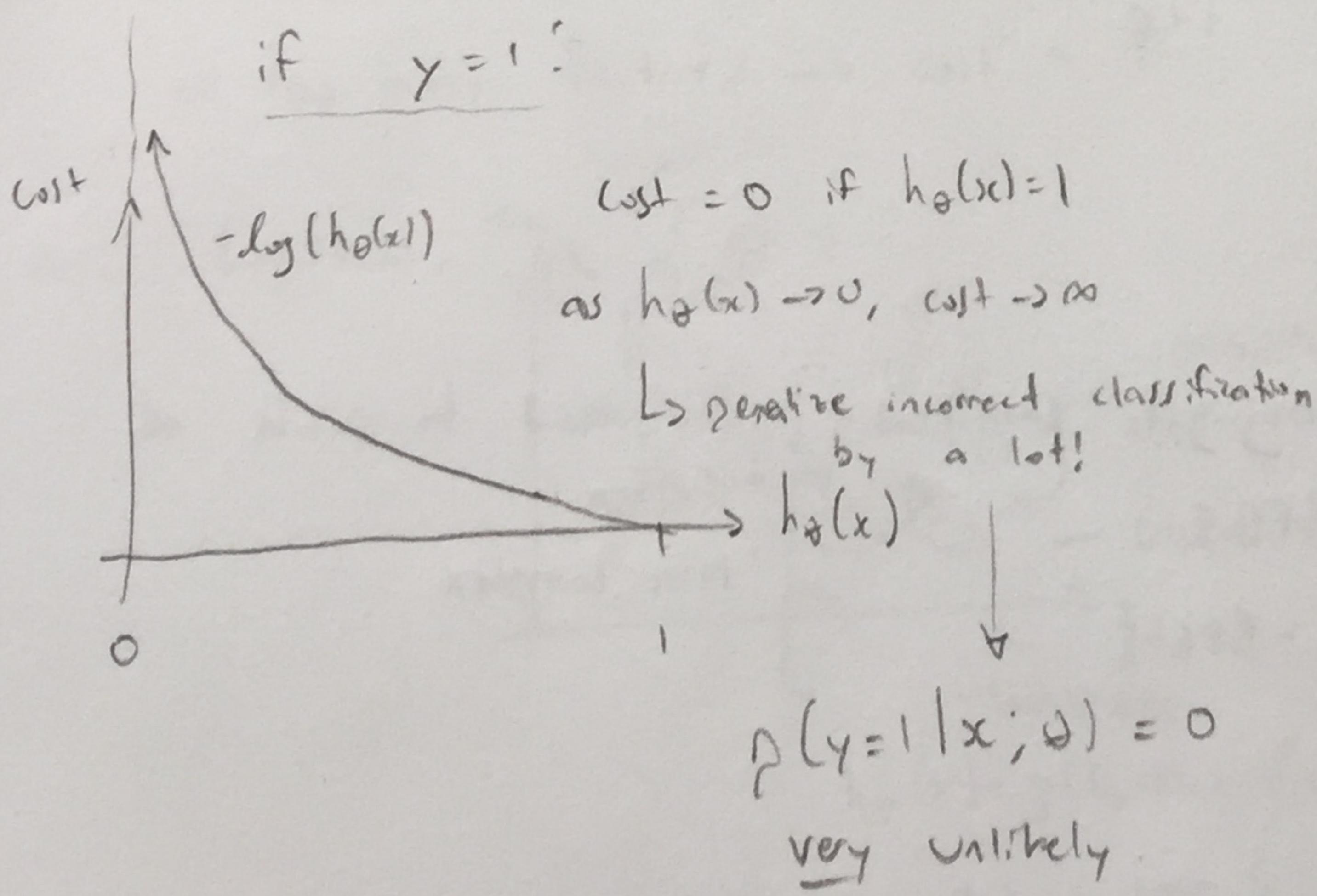
↳ for logistical, this results in a non-convex result



\Rightarrow b/c sigmoid $\frac{1}{1 + e^{-x}}$ is
nonlinear \rightarrow non-convex

\Rightarrow COST FUNCTION FOR LOGISTICAL REG:

$$\text{cost}(h_{\theta}(x), y) = \begin{cases} -\log h_{\theta}(x) & \text{if } y=1 \\ -\log(1-h_{\theta}(x)) & \text{if } y=0 \end{cases} \quad \text{GIVES CONVEX COST FNT.}$$



Simplified Cost Function & Gradient Descent

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(h_\theta(x^{(i)}), y^{(i)}) \quad \text{where} \quad \begin{cases} \text{cost} = -\log(h_\theta(x)) \text{ if } y=1 \\ \text{cost} = -\log(1-h_\theta(x)) \text{ if } y=0 \end{cases}$$

Note: $y=0$ or 1 always.

RE-WRITE: $\text{cost}(h_\theta(x), y) = -y \log(h_\theta(x)) - (1-y) \log(1-h_\theta(x))$

b/c y
always 0
or 1 .

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

↳ convex & derived from stats (principle of max. likel. est.)

↳ find θ that minimizes $J(\theta)$

$$\Rightarrow \text{to make prediction given new } x \rightarrow h_\theta(x) = \frac{1}{1+e^{-\theta^T x}} \quad (\gamma(y=1|x; \theta))$$

Gradient Descent:

$$\text{repeat } \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

Simultaneous update!

$$\Rightarrow \frac{\partial}{\partial \theta_j} J(\theta) = \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad \text{as before}$$

$$\text{here } h_\theta(x^{(i)}) = \frac{1}{1+e^{-\theta^T x}}$$

* Feature scaling also applies here.

Advanced Optimization

Other minimization of $J(\theta)$ algorithms:

- Conjugate Gradient
- BFGS
- L-BFGS

- No need to pick α
- Faster
- More complex

ex.) $\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \quad J(\theta) = (\theta_1 - 5)^2 + (\theta_2 - 5)^2$

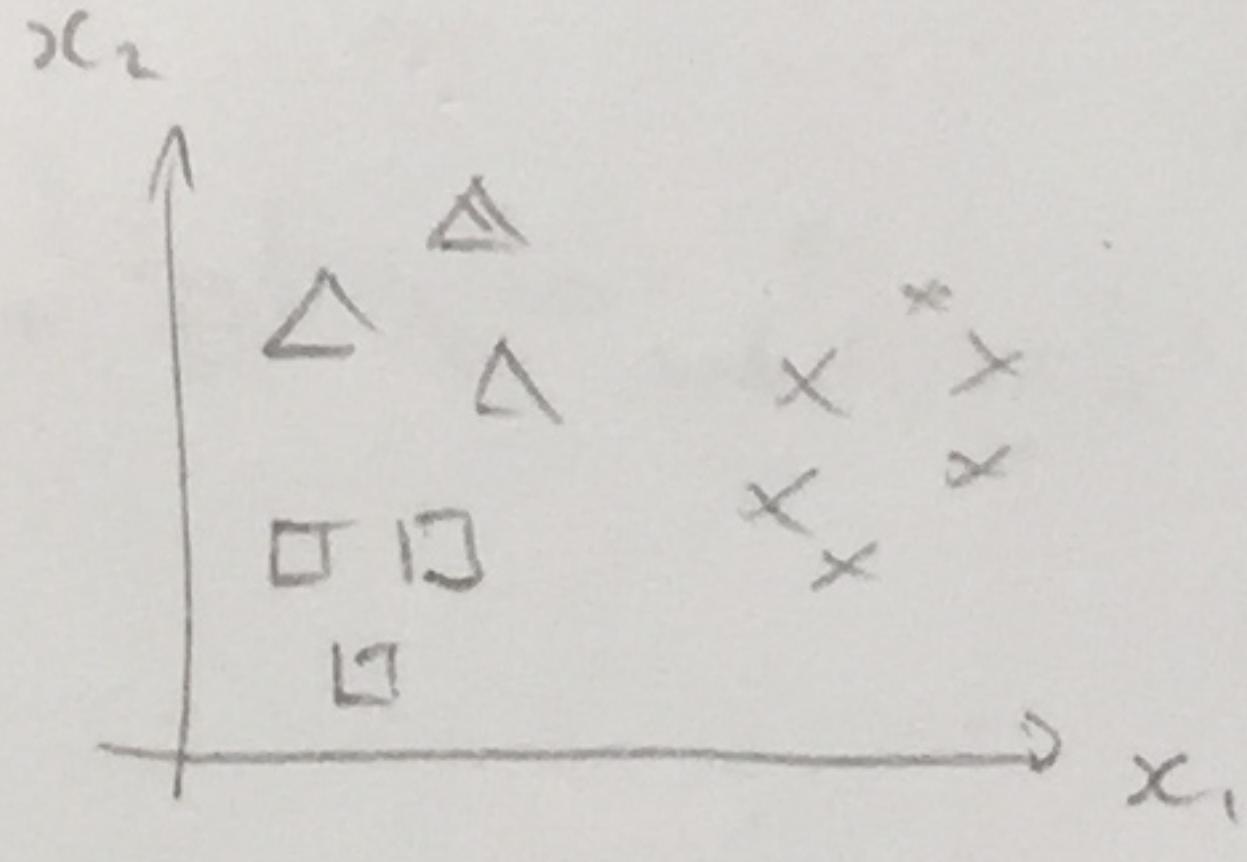
$$\frac{\partial}{\partial \theta_1} J(\theta) = 2(\theta_1 - 5) \quad \frac{\partial}{\partial \theta_2} J(\theta) = 2(\theta_2 - 5)$$

↳ see video for Octave implementation

Multi-class Classification:

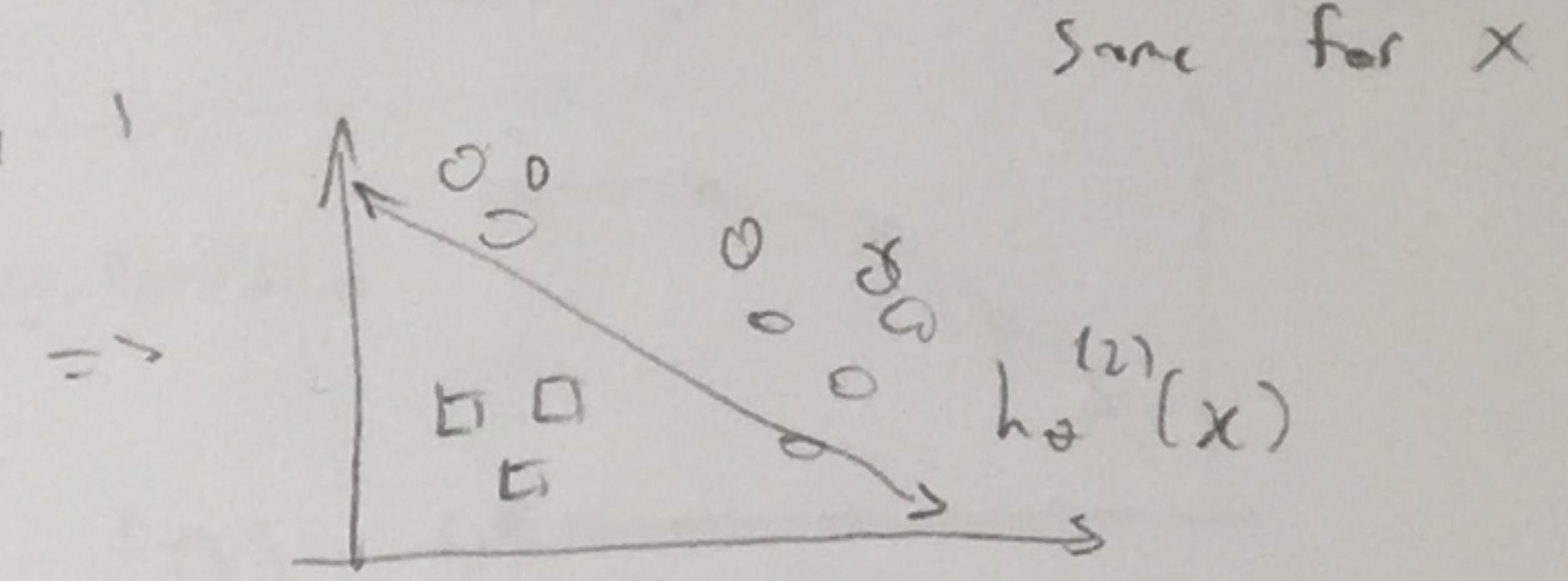
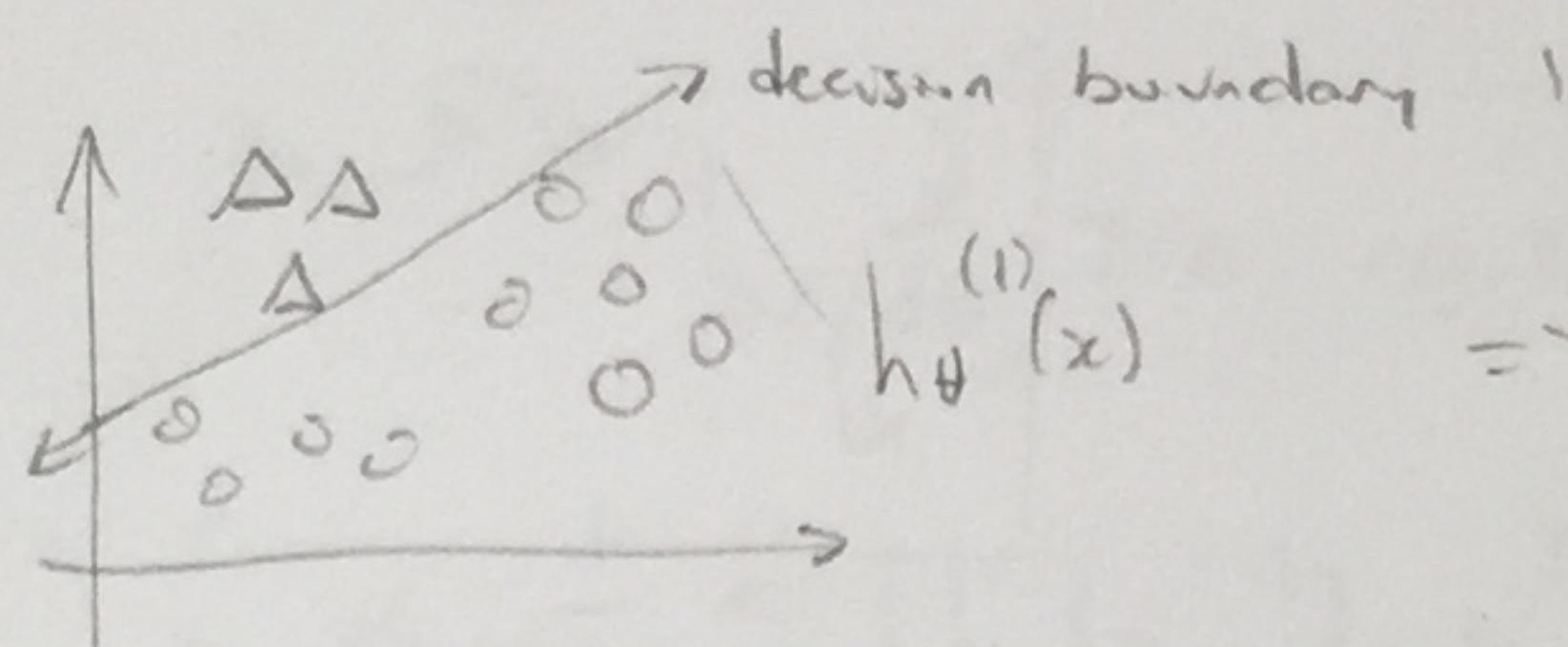
e.g.) email $\xrightarrow{\text{work}} \text{work } (1)$
 $\xrightarrow{\text{friends}} \text{friends } (2)$
 $\xrightarrow{\text{family}} \text{family } (3)$

not all
cold
flu
etc.



\Rightarrow One-vs-all (one vs rest)

Class 1) Δ :



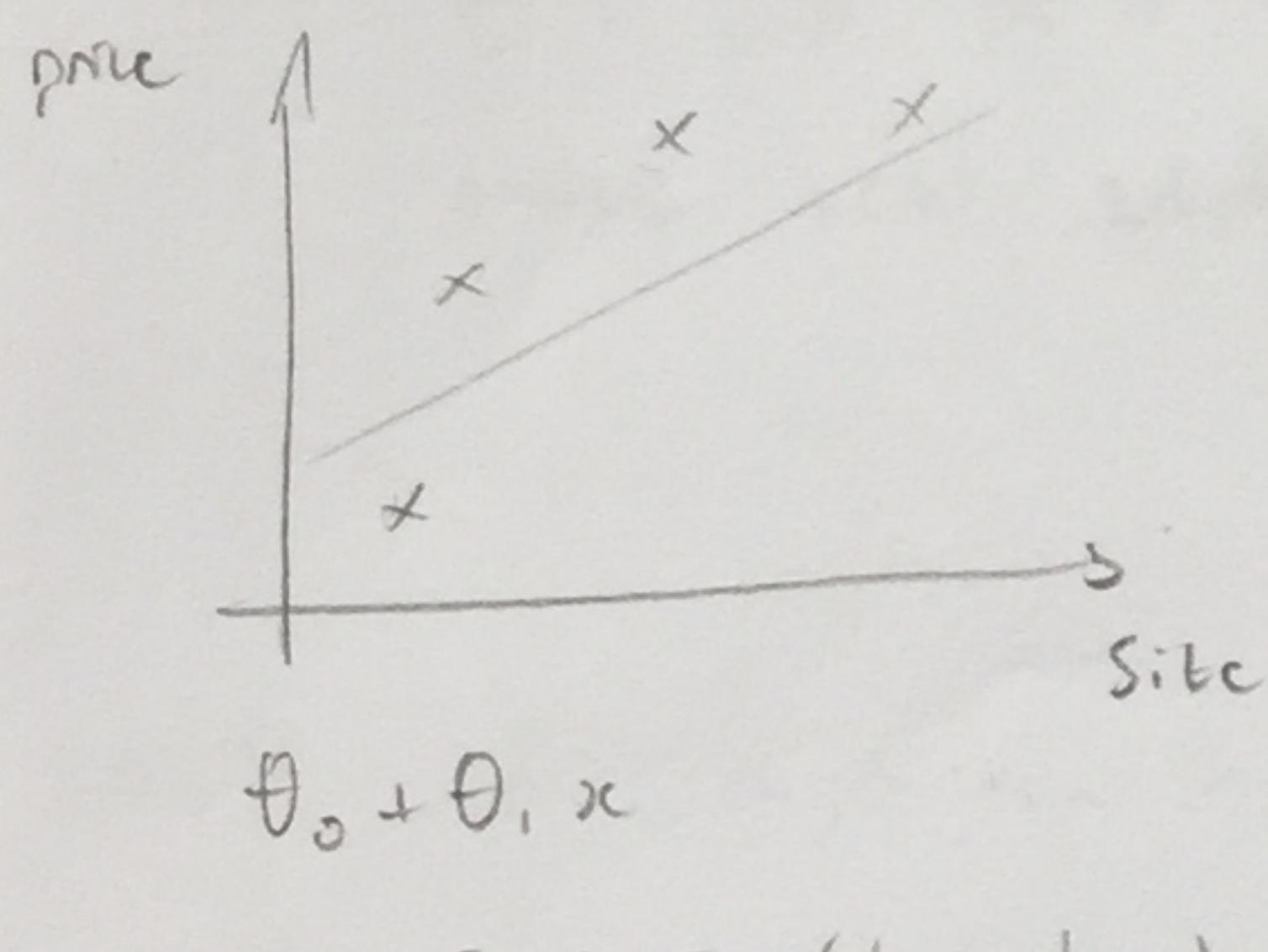
same for x

$$\Rightarrow h_{\theta}^{(i)}(x) = g(y=i | x; \theta), \quad i=1,2,3$$

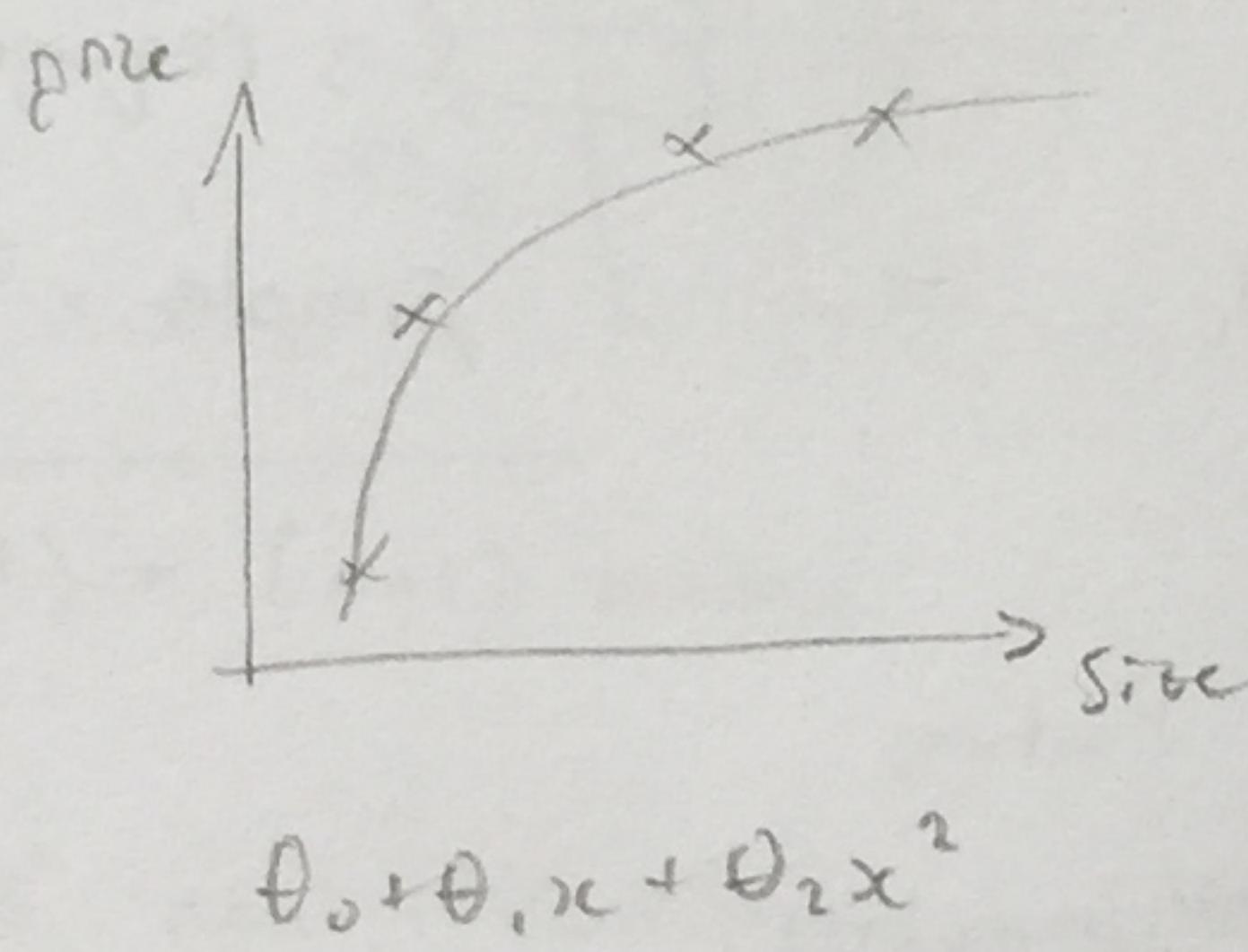
\hookrightarrow for new x , $\max_i h_{\theta}^{(i)}(x)$ gives prediction (e.g. 60%, 5%, 5%)

Don't add to 1 necessarily

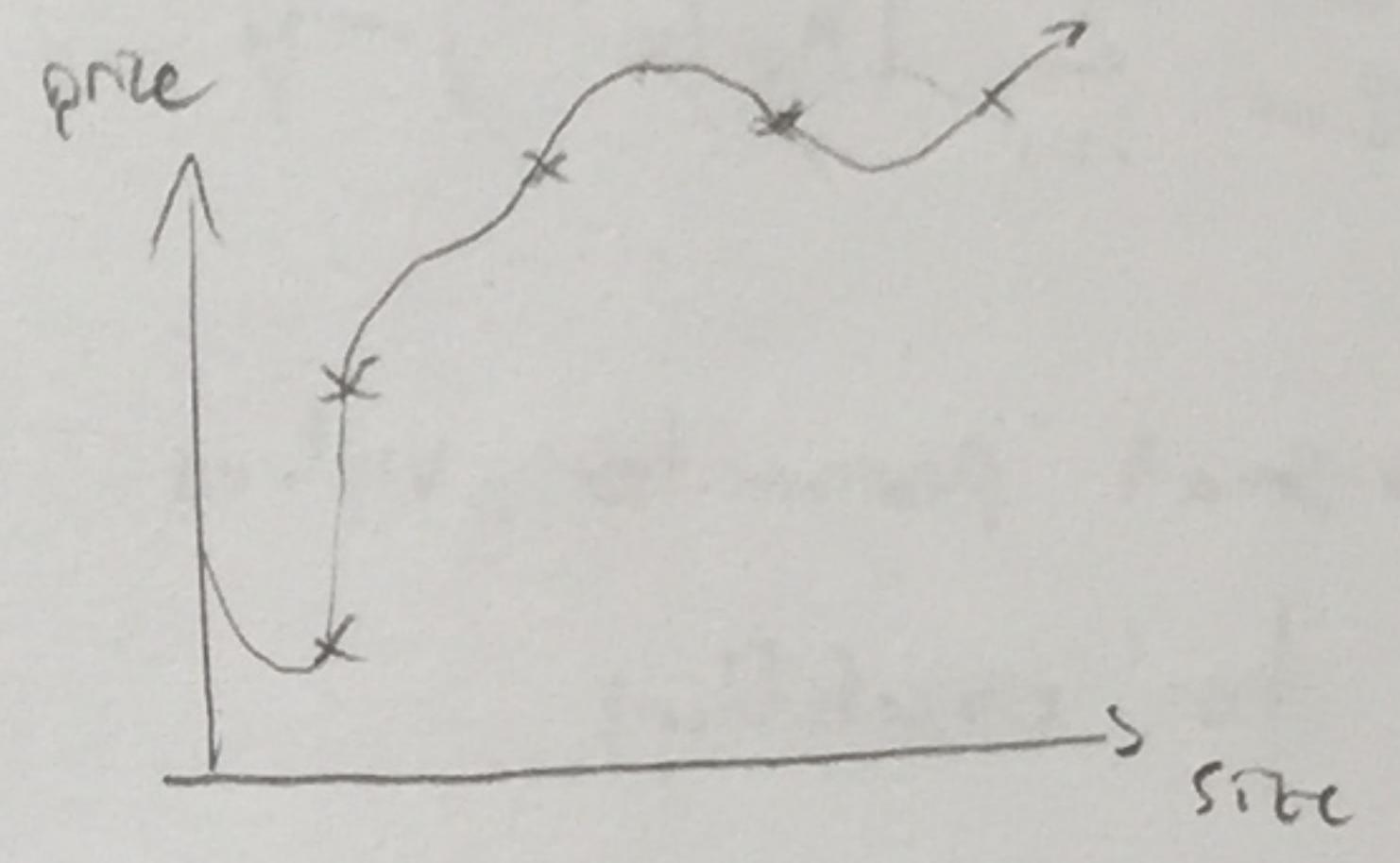
Overfitting:



UNDERFIT (high bias)



Pretty Good!

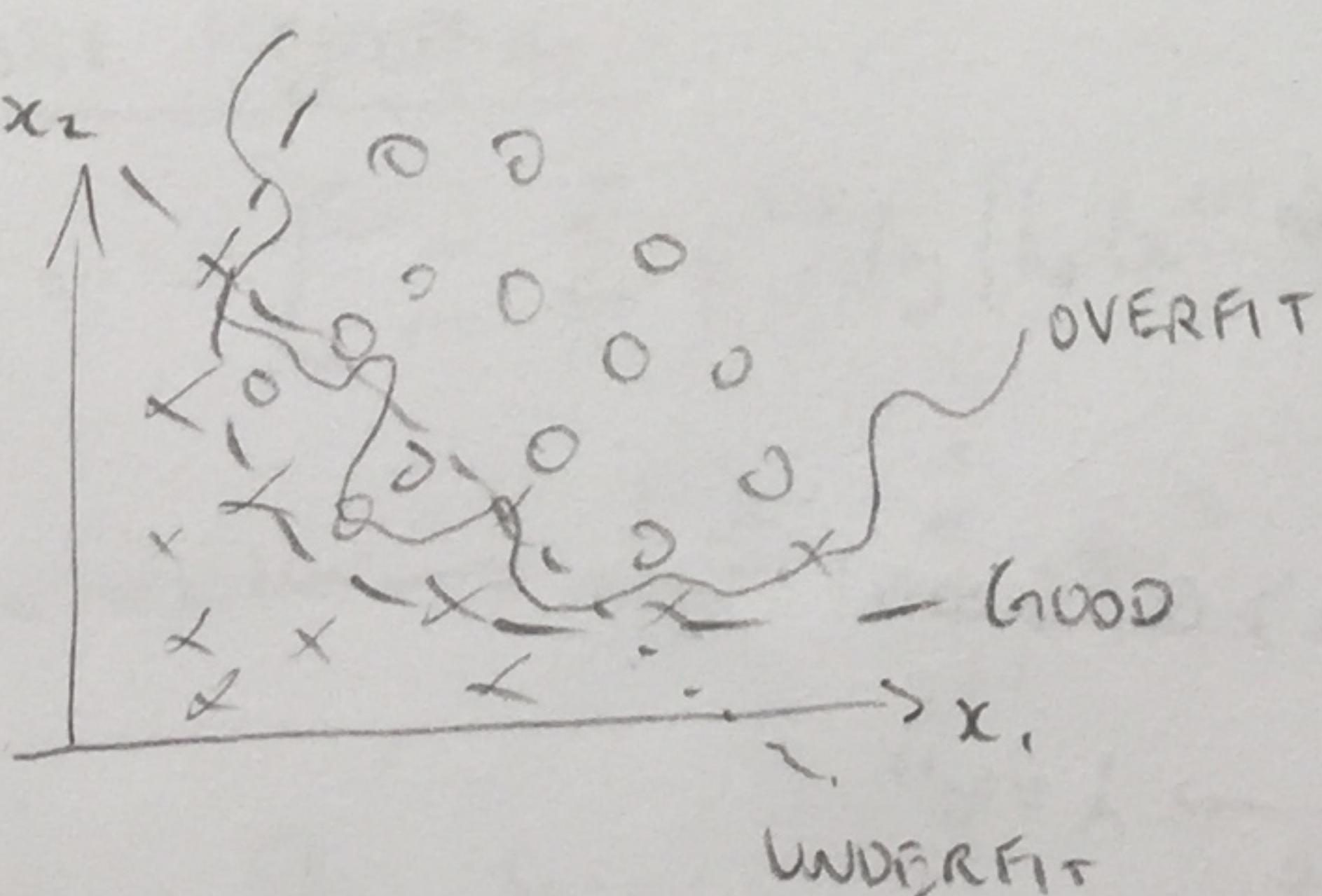


$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

OVERFIT (high variance)

\Rightarrow too many features \rightarrow cost ~ 0 , but fails to generalize to new examples!

ex2: Logistic:



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

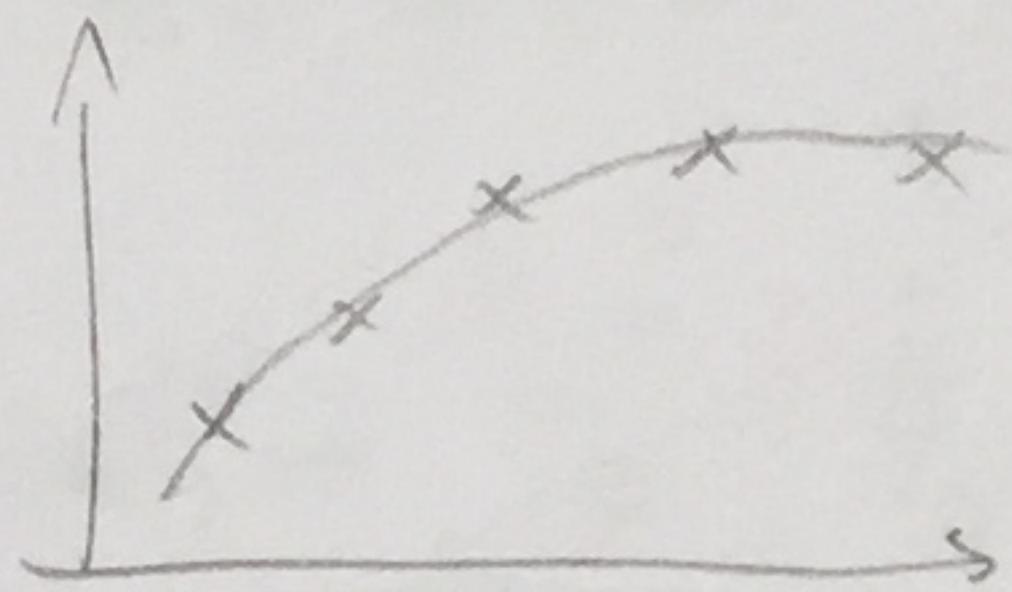
\Rightarrow we may just have a lot of features \rightarrow if we have too many & not a lot of training data, may overfit!

\hookrightarrow SOLTNS: \rightarrow reduce # features

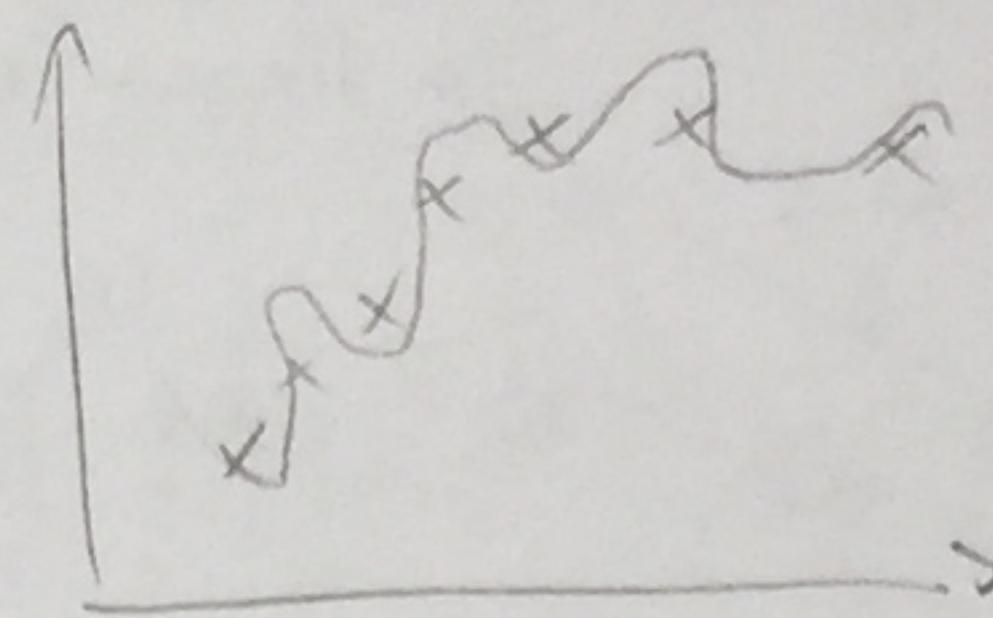
\rightarrow model select which features to keep / throw out

\rightarrow Regularization: keep all features, but change magnitudes of θ vals!

Regularization + Cost Function:



$$\theta_0 + \theta_1 x + \theta_2 x^2 \quad (\text{Quadratic})$$



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4 \rightarrow \text{Bad.}$$

\hookrightarrow penalize θ_3 & θ_4 , make v. small!

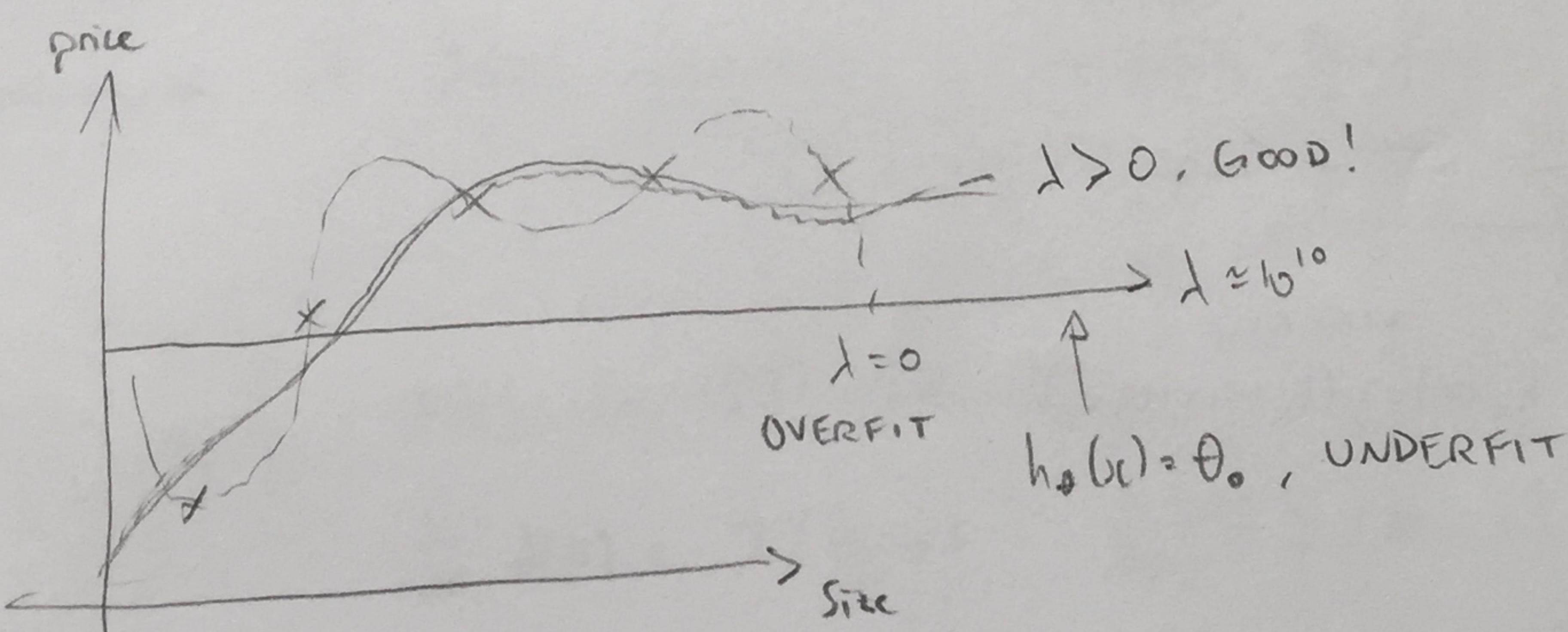
$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \underbrace{100\theta_3^2 + 100\theta_4^2}_{\text{require } \theta_3 \text{ & } \theta_4 \text{ to be small!}}$$

\hookrightarrow small parameter values $\theta_0, \theta_1, \dots, \theta_n$, simpler hypothesis less prone to overfitting

e.g. Housing: x_1, x_2, \dots, x_{100} Features
 $\theta_0, \theta_1, \theta_2, \dots, \theta_{100}$ Parameters

\hookrightarrow minimize all parameters: $J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2$
 \hookrightarrow not penalizing θ_0 by convention

$\hookrightarrow \lambda$: regularization parameter



Regularized Linear Regression:

New Cost Function: $J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$, goal $\min_{\theta} J(\theta)$

↳ Gradient Descent: $\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)} \leftarrow$ don't regularize θ_0 .

$$\theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \quad j=1, 2, \dots, n$$

$\frac{\partial}{\partial \theta_j} J(\theta)$ regularized

$$\rightarrow \theta_j := \underbrace{\theta_j \left(1 - \alpha \frac{\lambda}{m} \right)}_{1 - \alpha \frac{\lambda}{m} < 1} - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Normal Equation for Regularization:

$$\theta = (X^T X + \lambda \underbrace{\begin{bmatrix} 0 & & & & \\ & 1 & & & \\ & & 1 & & \\ & & & \ddots & \\ & & & & 0 \end{bmatrix}}_{(n+1) \times (n+1) \text{ matrix}})^{-1} X^T y$$

recall if $m \leq n \rightarrow X^T X$ non-invertible (degenerate, singular)

↳ luckily, if $\lambda > 0$, $X^T X + \lambda \begin{bmatrix} \dots \end{bmatrix} \rightarrow$ always invertible!

Regularized Logistic Regression:

$$J(\theta)_{\text{log.}} = - \left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)})) \right] + \underbrace{\frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2}_{\theta_1, \theta_2, \dots, \theta_n}$$

Gradient Descent: θ_0 same

$$\theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right], \quad j=1, 2, \dots, n$$

$$\hookrightarrow h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$