

(13)

WEEK

5

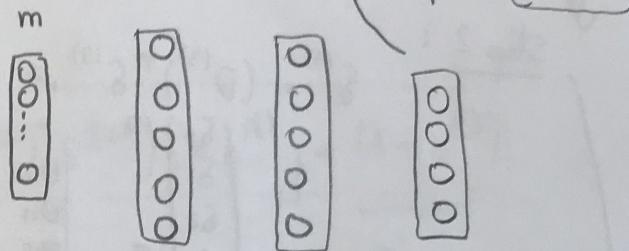
COST FUNCTION FOR NNs

↳ We sum over all output units (instead of only one for LG)

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K \left[y_k^{(i)} \log((h_\theta(x^{(i)}))_k) + (1 - y_k^{(i)}) \log(1 - (h_\theta(x^{(i)}))_k) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\theta_{j,i}^{(l)})^2$$

↙ Over all output units ↴ Reg. term.

$$\begin{cases} L \rightarrow \# \text{ of layers} \\ s_l \rightarrow \# \text{ of units in layer } l \\ K \rightarrow \# \text{ of output units} \end{cases}$$



Reg. term :

→ We do not count the bias units

→ In our example: $L=4$

$$\text{For } \theta^{(1)}: \rightarrow \sum_{i=1}^m \sum_{j=1}^5 (\theta_{j,i}^{(1)})^2$$

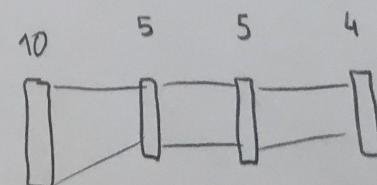
$$\theta^{(1)} \rightarrow \begin{bmatrix} \theta_{11}^{(1)} & \dots & \theta_{5m}^{(1)} \end{bmatrix} \quad \Bigg| \quad 5$$

$$\begin{cases} K = 4 \\ L = 4 \end{cases} \quad \begin{cases} s_1 = m \\ s_2 = 5 \\ s_3 = 5 \\ s_4 = K = 4 \end{cases}$$

$$\theta_1^{(1)} = \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \dots + \theta_{1m}^{(1)} x_m$$

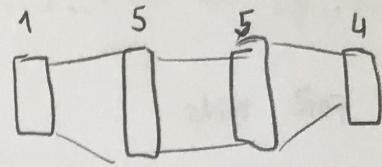
How many θ elements to find in the following NN?

$$\left. \begin{array}{l} \theta^{(1)} \rightarrow (10+1) \cdot 5 = 11 \cdot 5 = 55 \\ \theta^{(2)} \rightarrow (5+1) \cdot 5 = 6 \cdot 5 = 30 \\ \theta^{(3)} \rightarrow (5+1) \cdot 4 = 6 \cdot 4 = 24 \end{array} \right\} \underline{109} \text{ weights / params}$$



(14)

BACKPROP ALGO. FOR NNs

 (x, y) We need: $\rightarrow J(\theta)$

$$\rightarrow \frac{\partial}{\partial \theta_{ij}} J(\theta)$$

For each node we need to compute $\delta_j^{(l)}$ \rightarrow error of node j in layer l

Step 1: $\delta_j^{(4)} = a_j^{(4)} - y_j = h_\theta(x)_j - y_j$

$$\hookrightarrow \delta^{(4)} = a^{(4)} - y \rightarrow \begin{bmatrix} \delta_1 \\ \delta_2 \\ \delta_3 \\ \delta_4 \end{bmatrix} = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \end{bmatrix} - \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix}$$

Step 2: $\delta^{(3)} = (\theta^{(3)})^T \delta^{(4)} * g'(z^{(3)})$ $\theta^{(3)} \rightarrow 4 \times 6$

$$\hookrightarrow \begin{bmatrix} \delta_0 \\ \delta_1 \\ \delta_2 \\ \vdots \\ \delta_5 \end{bmatrix}^{(3)} = \begin{bmatrix} \theta_{10} & \dots & \theta_{14} \\ \theta_{11} & \ddots & \vdots \\ \theta_{12} & & \vdots \\ \vdots & & \vdots \\ \theta_{15} & \dots & \theta_{45} \end{bmatrix}^{(3)} \begin{bmatrix} \delta_1 \\ \delta_2 \\ \delta_3 \\ \delta_4 \end{bmatrix}^{(4)} * \begin{bmatrix} z_0 \\ z_1 \\ \vdots \\ z_5 \end{bmatrix}^{(3)}$$

$6 \times 1 \quad 6 \times 4 \quad 4 \times 1 \quad 6 \times 1$

Step 3: $\delta^{(2)} = (\theta^{(2)})^T \delta^{(3)} * g'(z^{(2)})$

(15)

ALGO IN DETAIL:

training set: $\{(x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})\}$ Set $\Delta_{ij}^{(l)} = 0$ for all i, j, l Loop over training samples $t = 1:m$ ① Set $a^{(1)} = x^{(t)}$ ② Perform forward prop. to get $a^{(l)}$ for $l = 2:L$

$$\begin{cases} z^{(2)} = \theta^{(1)} a^{(1)} \rightarrow a^{(2)} = g(z^{(2)}) \\ z^{(3)} = \theta^{(2)} a^{(2)} \rightarrow a^{(3)} = g(z^{(3)}) \\ z^{(4)} = \theta^{(3)} a^{(3)} \rightarrow a^{(4)} = h_{\theta}(x) = g(z^{(4)}) \end{cases}$$

③ Using $y^{(t)}$ → compute $\delta^{(L)} = a^{(L)} - y^{(t)}$ ④ Compute $\delta^{(L-1)} \rightarrow \delta^{(2)}$ using $\delta^{(l)} = ((\theta^{(l)})^T \delta^{(l+1)}) \cdot \underbrace{* a^{(l)} * (1 - a^{(l)})}_{g'}$

$$g'(z^{(1)}) = a^{(1)} \cdot (1 - a^{(1)})$$

$$⑤ \Delta_{ij}^{(l)} = \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$$

vect.

$$\Delta^{(l)} = \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$$

accumulator
of our partial.
der. terms

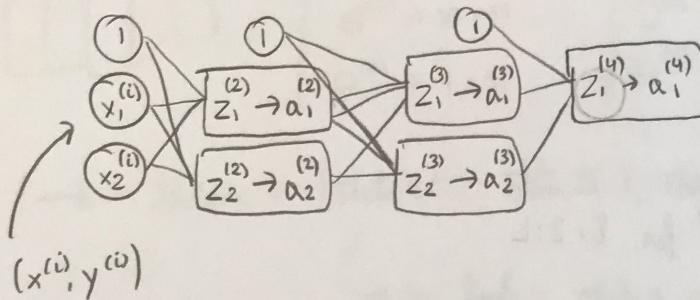
$$\left\{ \begin{array}{l} j \neq 0 \rightarrow D_{ij}^{(l)} = \frac{1}{m} (\Delta_{ij}^{(l)} + \lambda \theta_{ij}^{(l)}) \\ j = 0 \rightarrow D_{ij}^{(l)} = \frac{1}{m} \Delta_{ij}^{(l)} \end{array} \right.$$

$$\rightarrow \frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = D_{ij}^{(l)}$$

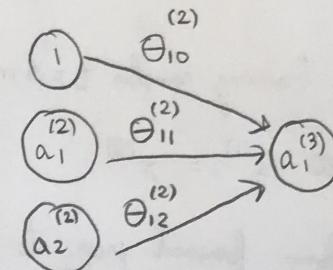
16

BACKPROP. INTUITION

→ Forward prop.



Weights:



$$z_1^{(3)} = \theta_{10}^{(2)} \cdot 1 + \theta_{11}^{(2)} \cdot a_1^{(2)} + \theta_{12}^{(2)} \cdot a_2^{(2)}$$

→ Back. prop.

Special case

- one example $x^{(i)}, y^{(i)}$
- 1 output unit
- $\lambda = 0$

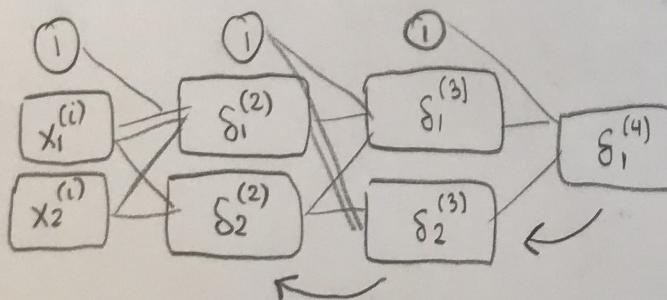
$$\text{cost}(i) = y^{(i)} \log \sigma(x^{(i)}) + (1 - y^{(i)}) \log(1 - \sigma(x^{(i)}))$$

↳ how well is it doing on example (i) ?

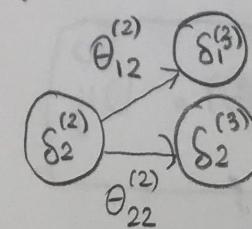
→ $\delta_j^{(l)}$ → error of cost for $a_j^{(l)}$ (unit j in layer l)

$$\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(i)$$

→ measure of how much we would like to change the weights of the NN to affect these values which will affect the output and thus the cost.



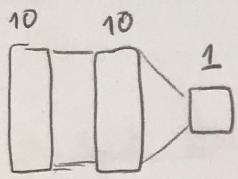
Weights:



$$\delta_2^{(2)} = \theta_{12}^{(2)} \delta_1^{(3)} + \theta_{22}^{(2)} \delta_2^{(3)}$$

(17)

Unrolling the params (example):



$$\left\{ \begin{array}{l} \theta^{(1)} \rightarrow 10 \times 11 \\ \theta^{(2)} \rightarrow 10 \times 11 \\ \theta^{(3)} \rightarrow 1 \times 11 \end{array} \right.$$

$$\left\{ \begin{array}{l} D^{(1)} \rightarrow 10 \times 11 \\ D^{(2)} \rightarrow 10 \times 11 \\ D^{(3)} \rightarrow 1 \times 11 \end{array} \right.$$

$$\rightarrow \text{thetaVec} = [\text{Theta1}(:); \text{Theta2}(:); \text{Theta3}(:)];$$

↳ Go back: $\text{Theta1} = \text{reshape}(\text{thetaVec}(1:110), 10, 11);$

In learning algo:

→ Have initial $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$

→ Unroll for fminunc (λ costFunction, initialTheta, options)

→ gradientVec \leftarrow costFunction(thetaVec)

→ Get $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$ from thetaVec (reshape)

→ Compute $D^{(1)}, D^{(2)}, D^{(3)}$

→ Unroll Ds to get gradientVec

Note:

→ Matrix form: better for FP and BP computations

→ Vector form: for functions that assume you have all your params. in one input

(18)

GRADIENT CHECKING

↳ to avoid implementation bugs

Estimation : $\frac{d}{d\theta} J(\theta) \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$ $\epsilon = 10^{-4}$

For θ vector with multiple params:

$$\frac{\partial}{\partial \theta_2} J(\theta) \approx \frac{J(\theta_1, \theta_2 + \epsilon, \dots, \theta_n) - J(\theta_1, \theta_2 - \epsilon, \dots, \theta_n)}{2\epsilon}$$

↳ Implement estimation and compare with Ds from backprop.

Random init.

- Do not init. your $\theta_{ij}^{(l)}$ to 0s $\rightarrow \alpha_1^{(2)} = \alpha_2^{(2)}$ \rightarrow same features learned everywhere
- Pick random values $[-\epsilon, \epsilon]$ \rightarrow symmetry breaking

OVERALL METHOD:

① Pick a NN architecture

Training :

- ① Weights init. randomly
- ② Implement forward prop. \rightarrow get $h_\theta(x^{(i)})$ for all $x^{(i)}$
- ③ Implement cost function
- ④ Implement BP \rightarrow get partial derivatives
- ⑤ Use grad. checking to confirm BP
- ⑥ Use GD or other to minimise cost function \rightarrow get $\theta_{ij}^{(l)}$

(19)

WEEK

(6)

EVALUATING A HYPOTHESIS:

→ Split training/test sets

- ↳ Learn θ and minimize $J_{\text{train}}(\theta)$ with training set
- ↳ Compute the test set error $J_{\text{test}}(\theta)$

→ Test set error:

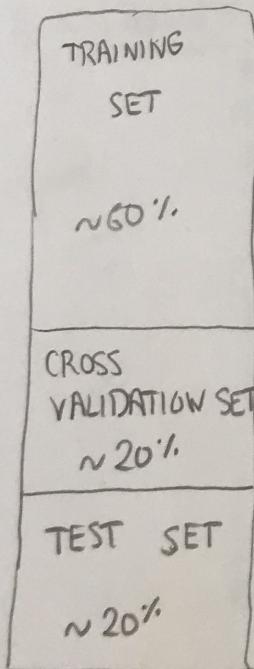
$$\textcircled{1} \text{ Lin. reg. } \rightarrow J_{\text{test}}(\theta) = \frac{1}{2m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} (h_{\theta}(x_{\text{test}}^{(i)}) - y_{\text{test}}^{(i)})^2$$

\textcircled{2} Classification → Misclassification error

$$\text{error}(h_{\theta}(x), y) = \begin{cases} 1 & \text{when } h_{\theta}(x) \geq 0.5 \text{ and } y=0 \text{ or } h_{\theta}(x) < 0.5 \text{ and } y=1 \\ 0 & \text{otherwise} \end{cases}$$

$$\text{Test Error} = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \text{err}(h_{\theta}(x_{\text{test}}^{(i)}), y_{\text{test}}^{(i)})$$

MODEL SELECTION:



\textcircled{1} Optimize params. in θ by using training set for each model.

\textcircled{2} Find the model with smaller error using the CV set.

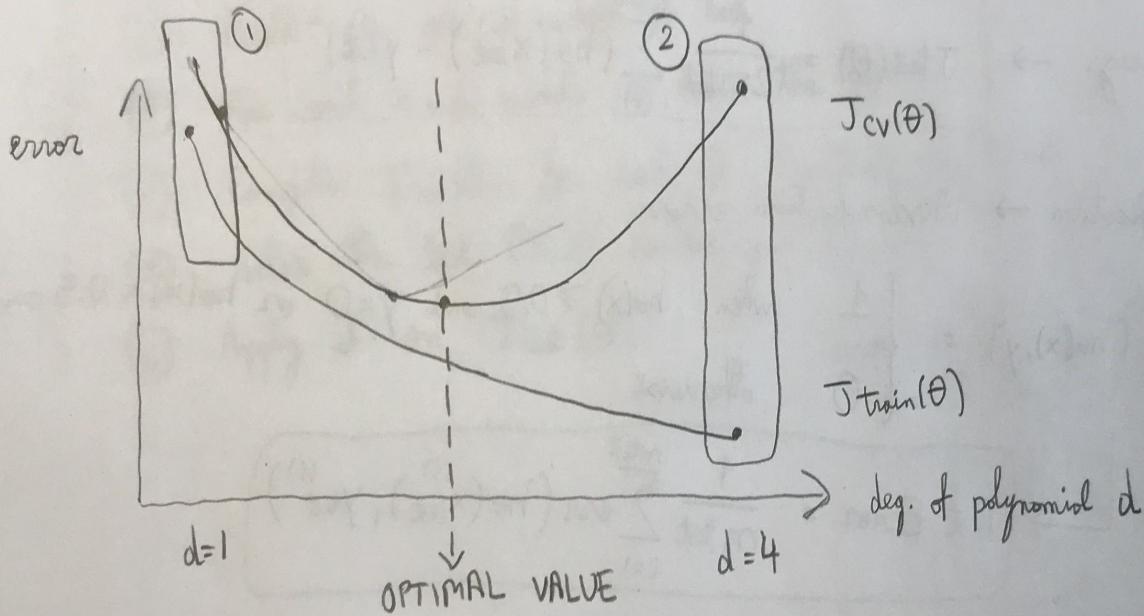
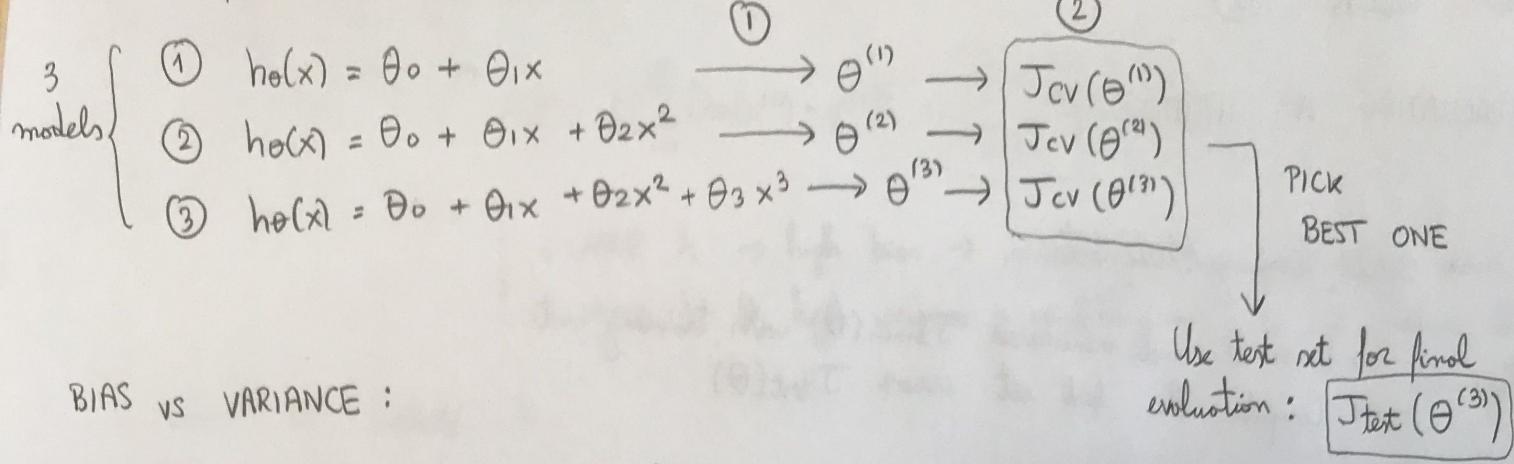
\textcircled{3} Estimate the generalization error using the test set

* Usually $J_{\text{cv}}(\theta) < J_{\text{test}}(\theta)$ because the degree of the polynomial has been fit to the CV set.

PROCEDURE

(20)

Example:



① HIGH BIAS → underfitting

- $J_{train}(\theta)$ will be high
- $J_{CV}(\theta) \approx J_{train}(\theta)$

② HIGH VARIANCE → overfitting

- $J_{train}(\theta)$ will be low
- $J_{CV}(\theta) \gg J_{train}(\theta)$

(21)

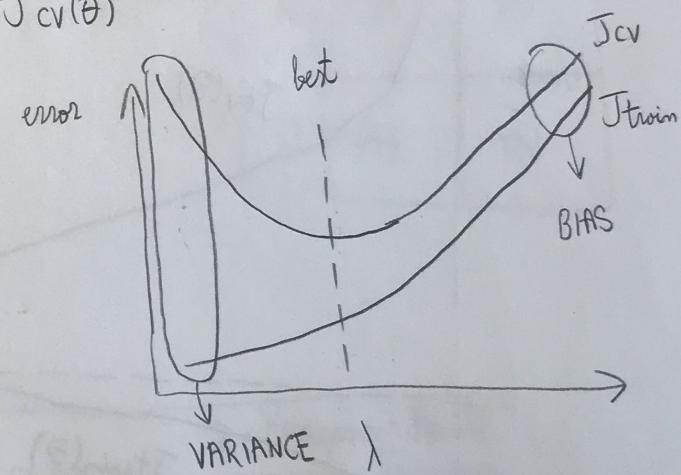
→ Selecting regularization param. λ :

$$\text{LIN. REG.} \rightarrow J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \boxed{\frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2}$$

| → LARGE $\lambda \rightarrow$ high bias \rightarrow underfitting \rightarrow we penalize too much
 | → SMALL $\lambda \rightarrow$ high variance \rightarrow overfitting $\rightarrow \lambda$ is not affecting the optimization

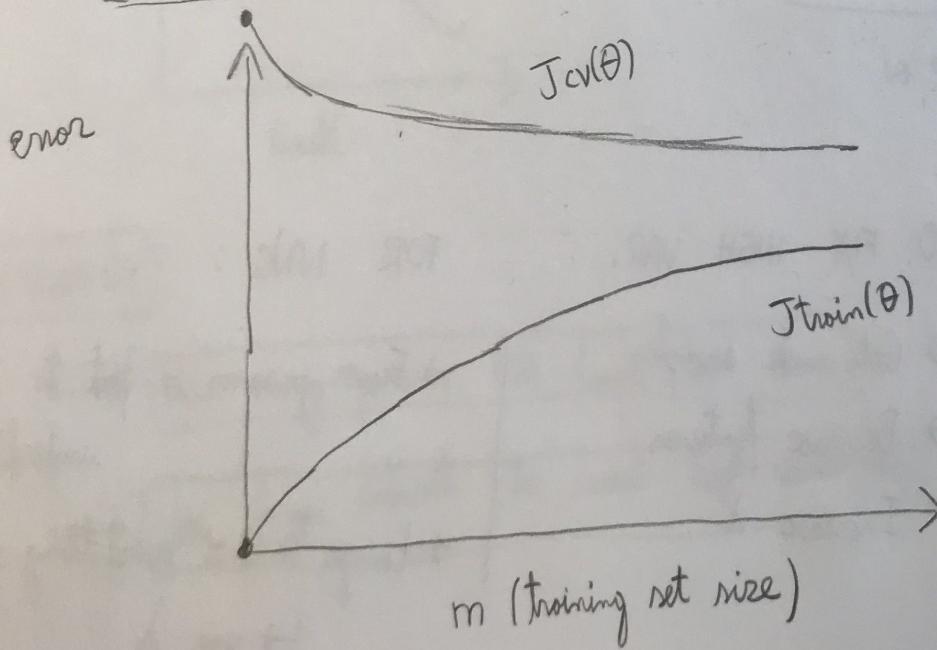
HOW TO CHOOSE λ :

- ① Create list of λ values $\lambda = 0, 0.01, 0.02, 0.04, \dots, 10$
- ② Create your set of models to test
- ③ For each λ and models \rightarrow learn the θ
- ④ Compute $J_{cv}(\theta)$ for each θ (without reg. term)
- ⑤ Select the best (θ, λ) combo on $J_{cv}(\theta)$
- ⑥ Apply it on $J_{test}(\theta)$



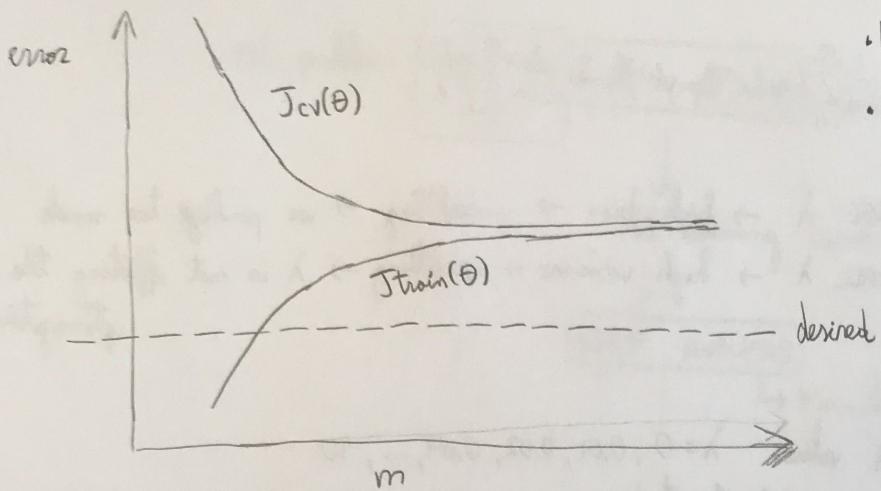
TYPICAL LEARNING CURVE EXPECTED:

$$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$



22

HIGH BIAS :

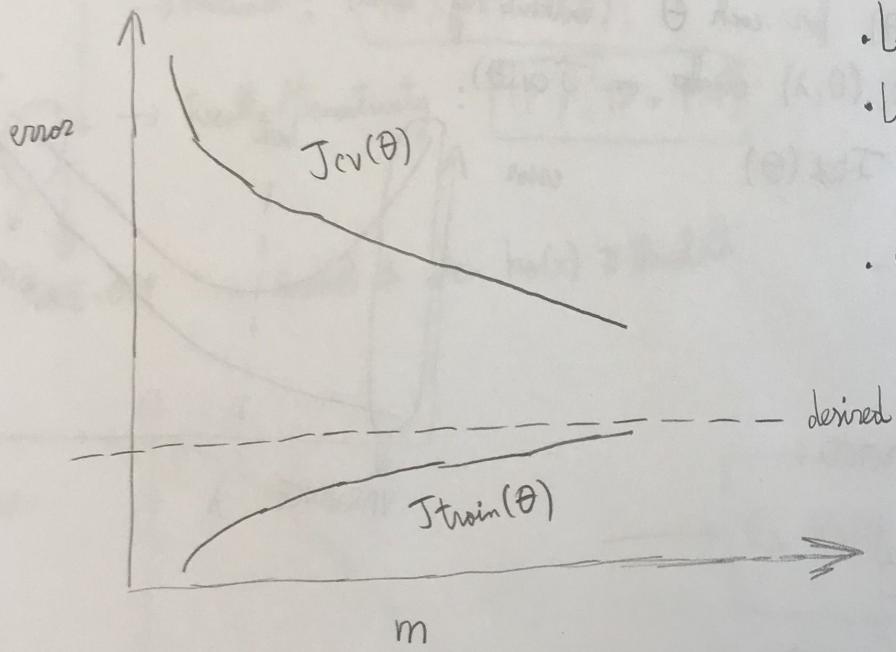


- Low $m \rightarrow J_{train}$ low and J_{cv} high
- Large $m \rightarrow$ Both J_{train} and J_{cv} high
and $J_{train} \approx J_{cv}$



Having larger dataset
WILL NOT HELP

HIGH VARIANCE :



- Low $m \rightarrow J_{train}$ low and J_{cv} high
- Large $m \rightarrow J_{train}$ gets higher while
 J_{cv} get lower and $J_{train} < J_{cv}$
- BUT $|J_{cv} - J_{train}|$ remains high



Having larger dataset
WILL PROBABLY HELP

TIPS

FIX HIGH BIAS

- Add features
- Add polynomial features
- Decrease λ

TO FIX HIGH VAR.

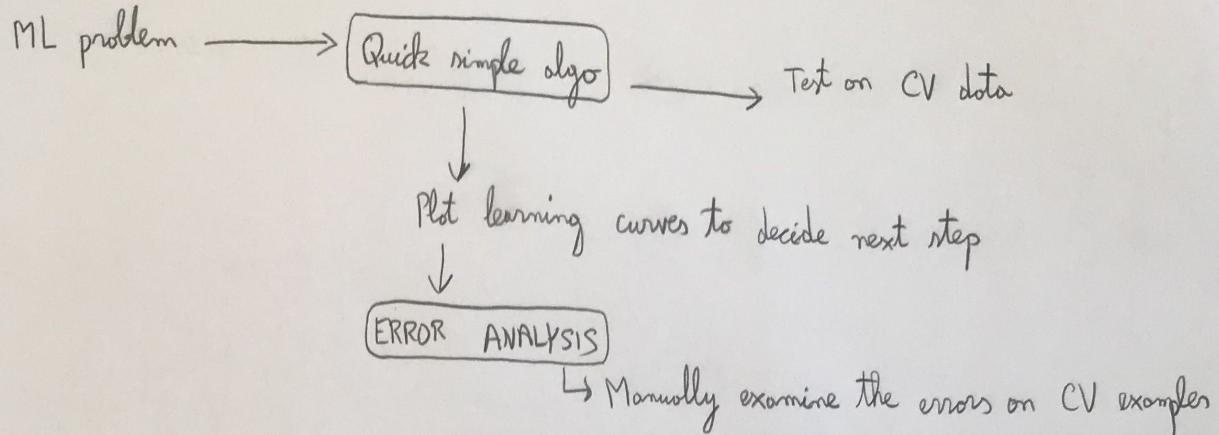
- Get more samples
- Decrease features
- Increase λ

FOR NNs :

- Fewer params. = fast to train but underfitting
- Larger NN = overfitting and slow
↳ use λ

(23)

ERROR ANALYSIS:

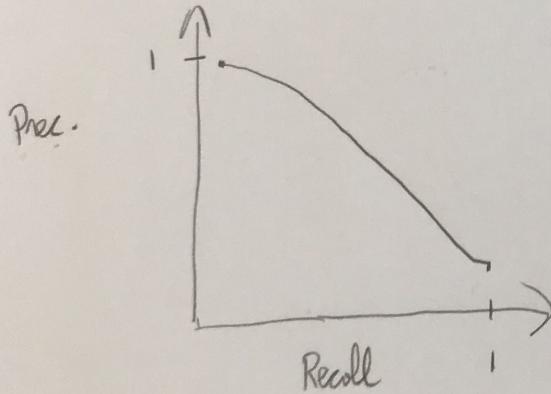


CLASS IMBALANCE (SKEWED CLASSES):

$$\left\{ \begin{array}{l} \rightarrow \text{Precision : } \frac{\text{TP}}{\text{TP} + \text{FP}} \\ \rightarrow \text{Recall / Sensitivity : } \frac{\text{TP}}{\text{TP} + \text{FN}} \end{array} \right.$$

↓
TRADE-OFF → Predict 1 when $h_0(x) \geq \text{threshold}$

		Ground Truth	
		1	0
Predicted	1	TP	FP
	0	FN	TN



$$F_1 \text{ SCORE : } \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

↳ Select this on the CV set!

STRATEGY

- ↳ Use many params. (to have low bias)
- ↳ Use huge dataset (to avoid high variance).