

What to Try Next?

e.g. implemented regularized lin. regression to predict housing prices - makes large errors on new set of houses! $(J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right])$

- ↳ POSSIBILITIES:
- get more training data
 - try smaller feature sets, additional features, polynomial features
 - decrease/increase λ (play w/ under/overfitting)

⇒ SIMPLE TECHNIQUES: ML DIAGNOSTICS

- ↳ gain insight into what is /isn't working
 ↳ time costly, but helpful

Evaluating a Hypothesis:

⇒ low training error \neq good $h_\theta(x)$! (overfitting)

⇒ split data into train + test (e.g. 70/30)

$$(x^{(i)}, y^{(i)}) \quad (x_{\text{test}}^{(i)}, y_{\text{test}}^{(i)}) \quad \Rightarrow \text{Send randomly, "hold out" method}$$

↓ ↓

⇒ process: ① get $h_\theta(x)$ from training set

$$\text{② compute test set error } J_{\text{test}}(\theta) = \frac{1}{2M_{\text{test}}} \sum_{i=1}^{M_{\text{test}}} (h_\theta(x_{\text{test}}^{(i)}) - y_{\text{test}}^{(i)})^2$$

↳ for logistic regression, $J_{\text{test}}(\theta) = -\frac{1}{M_{\text{test}}} \left[\sum_{i=1}^{M_{\text{test}}} y_{\text{test}}^{(i)} \log(h_\theta(x_{\text{test}}^{(i)})) + (1-y_{\text{test}}^{(i)}) \log(1-h_\theta(x_{\text{test}}^{(i)})) \right]$

↳ could also compute "misclassification error" for logistic reg.

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

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

Choosing λ : Model Selection Problems:

Choosing b/n: $h_{\theta}(x) = \theta_0 + \theta_1 x$ deg = 1
 $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$ deg = 2
 \vdots
 $h_{\theta}(x)$ deg = 10
 \Rightarrow what is best?

$\Rightarrow \Theta^{(1)}$ is $[2 \times 1]$ (degree = 1)
 $\Theta^{(2)}$ is $[3 \times 1]$ (degree = 2) \Rightarrow could compare $J(\theta^d)$'s!
 (test set error)

\Rightarrow BAD comparison \rightarrow NOT predictive for fitting to new models!

\Rightarrow Better Method: $\sim 60\%$ data = training set
 $\sim 20\%$ data = cross validation set
 $\sim 20\%$ data = test set] split data

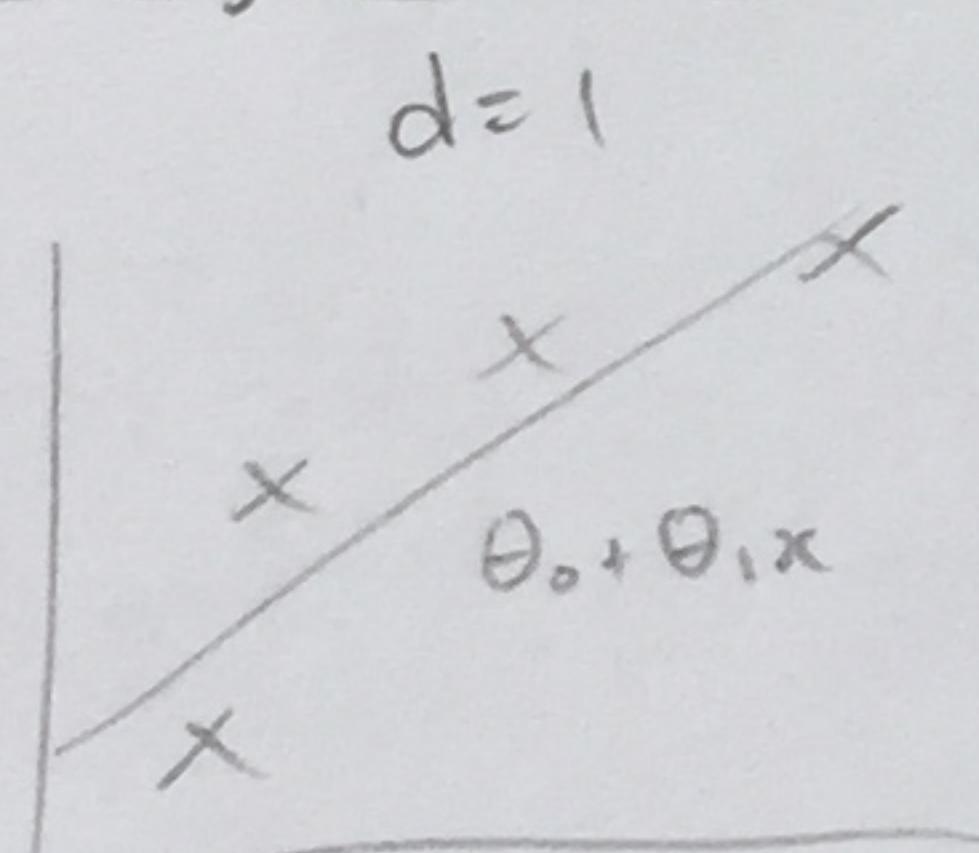
\Rightarrow Mcv examples of $(x_{cv}^{(i)}, y_{cv}^{(i)})$, etc.

$\Rightarrow J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$ & likewise w/ $J_{cv}(\theta)$, $J_{test}(\theta)$

Now can compare $J_{cv}(\theta^{(d)}) \rightarrow$ choose deg w/ lowest $J(\theta^{(d)})$

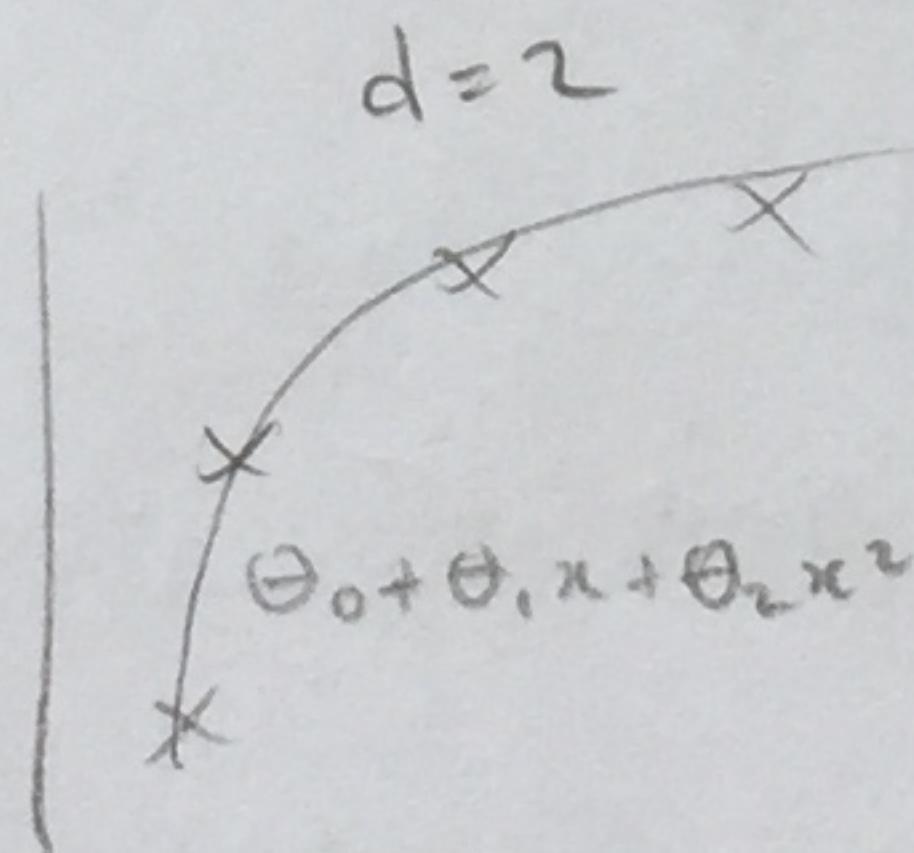
\hookrightarrow report out general error using $J_{test}(\theta^d)$ (test set)

Diagnosing Bias & Variance:

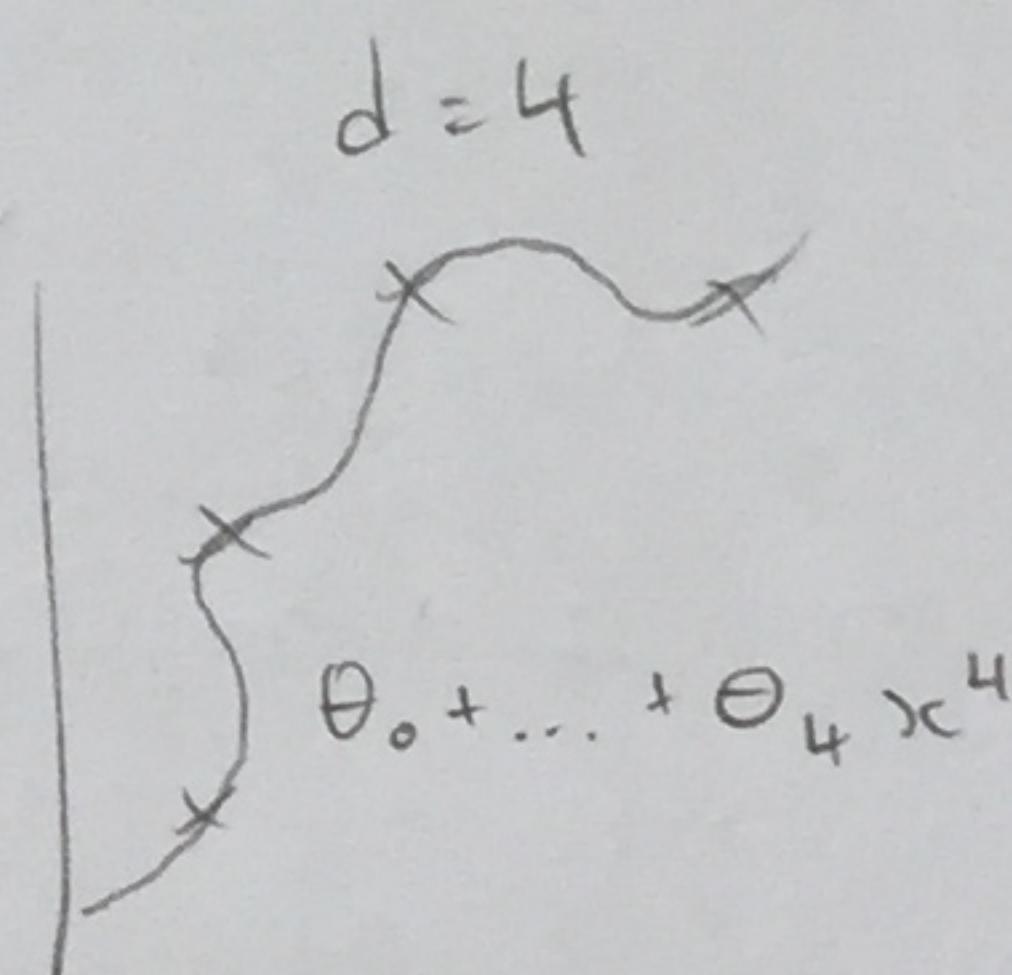


High Bias
(Underfit)

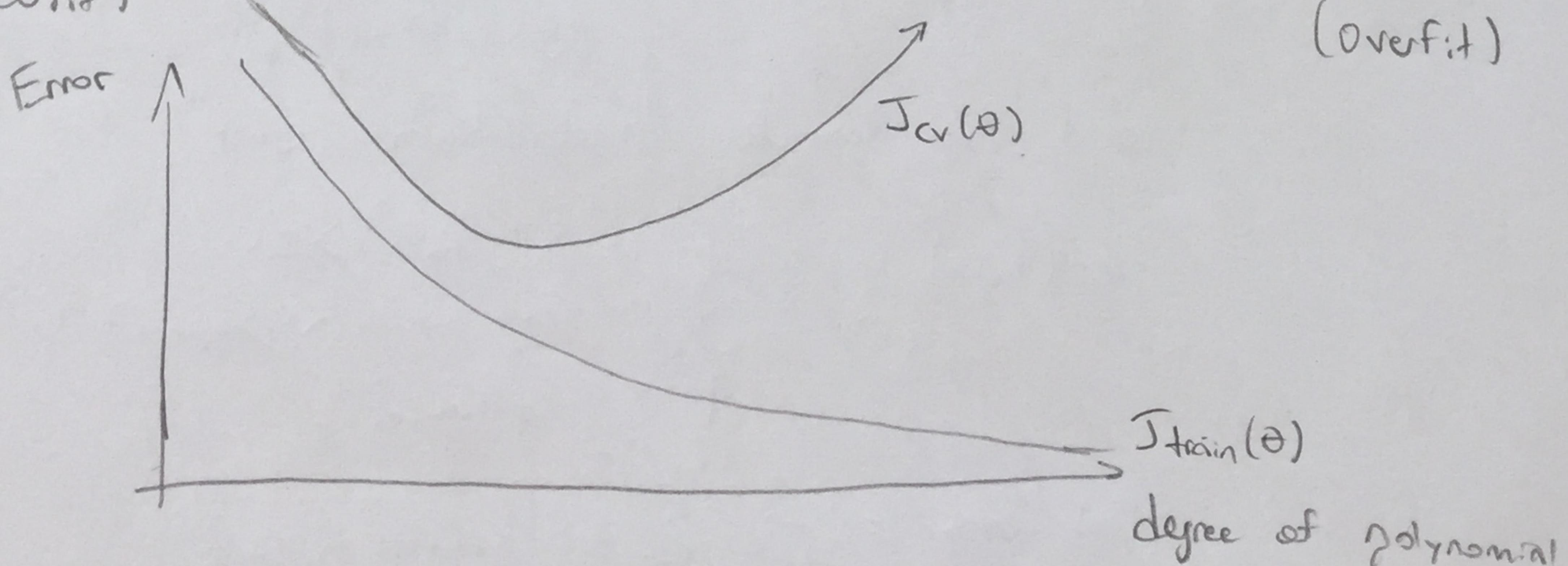
vs.



Just Right



High Variance
(Overfit)



\Rightarrow if $J_{cv}(\theta)$ is high, are we over (variance) or under (bias) fitting?

\hookrightarrow HIGH BIAS: J_{train} & J_{cv} both high ($J_{cv} \approx J_{train}$)

HIGH VARIANCE: J_{train} low but J_{cv} high ($J_{cv} \gg J_{train}$)

Regularization & Bias / Variance:

$$h_\theta(x) = \theta_0 + \dots + \theta_d x^d$$

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

$\hookrightarrow \uparrow \lambda \rightarrow$ high bias (underfit)

mid $\lambda \rightarrow$ just right

$\downarrow \lambda \rightarrow$ high variance (overfit)

Note $J_{train}, J_{cv}, J_{test}$ all same, just use $M_{train}, M_{cv}, M_{test}$

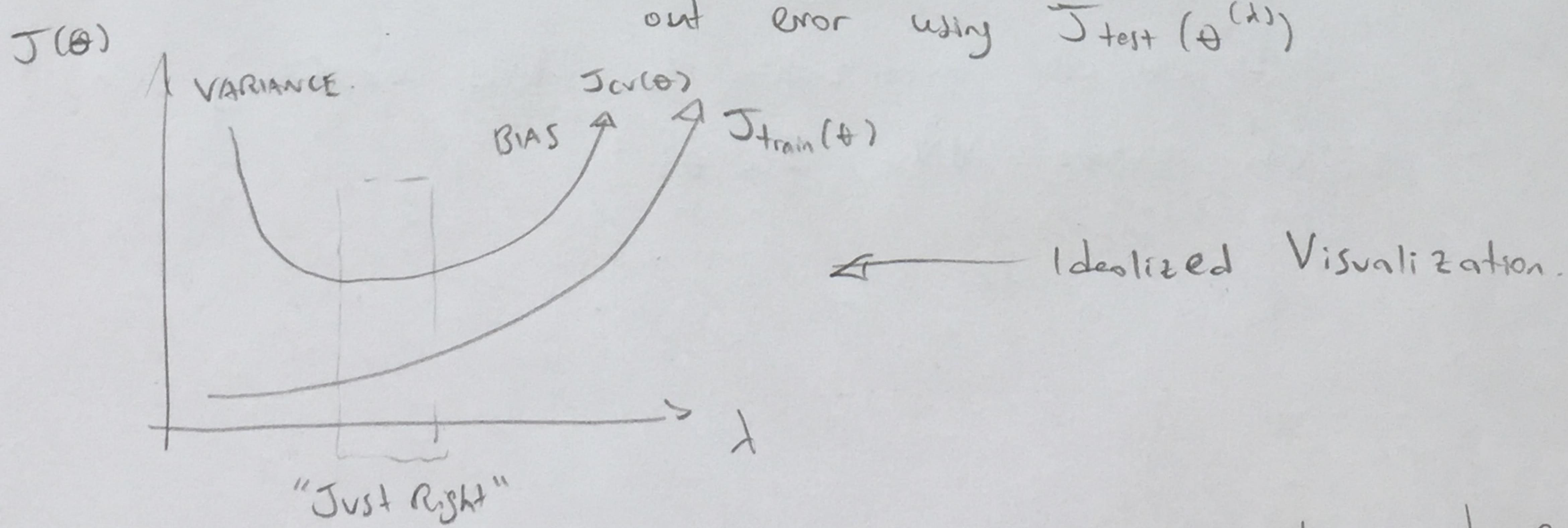
\Rightarrow to choose λ : cycle through λ values (eg. 0, 0.02, 0.04 \rightarrow 10.0)

$\lambda = 0 \rightarrow J(\theta)$, set $\theta^{(1)} \leftarrow$ # of λ , not degree

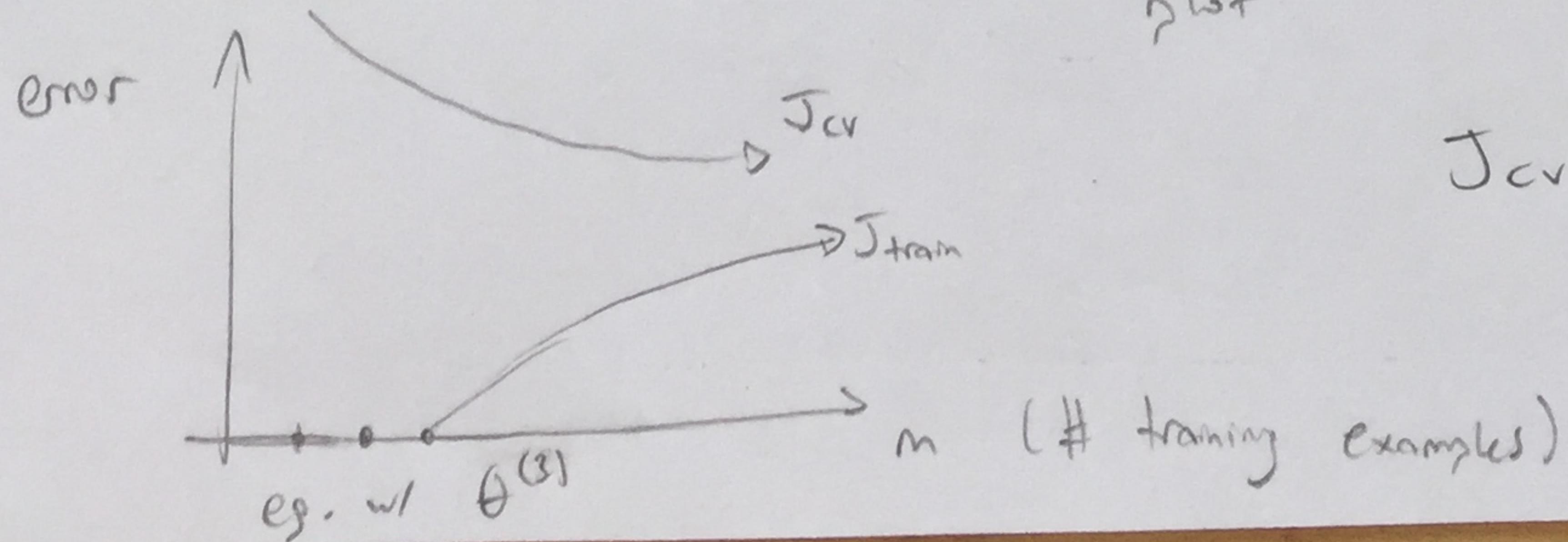
$\lambda = 0.02 \rightarrow J(\theta)$, get $\theta^{(2)}$

etc.

\hookrightarrow pick λ that gives lowest $J_{cv}(\theta^{(\lambda)})$, report out error using $J_{test}(\theta^{(\lambda)})$



Learning Curves:

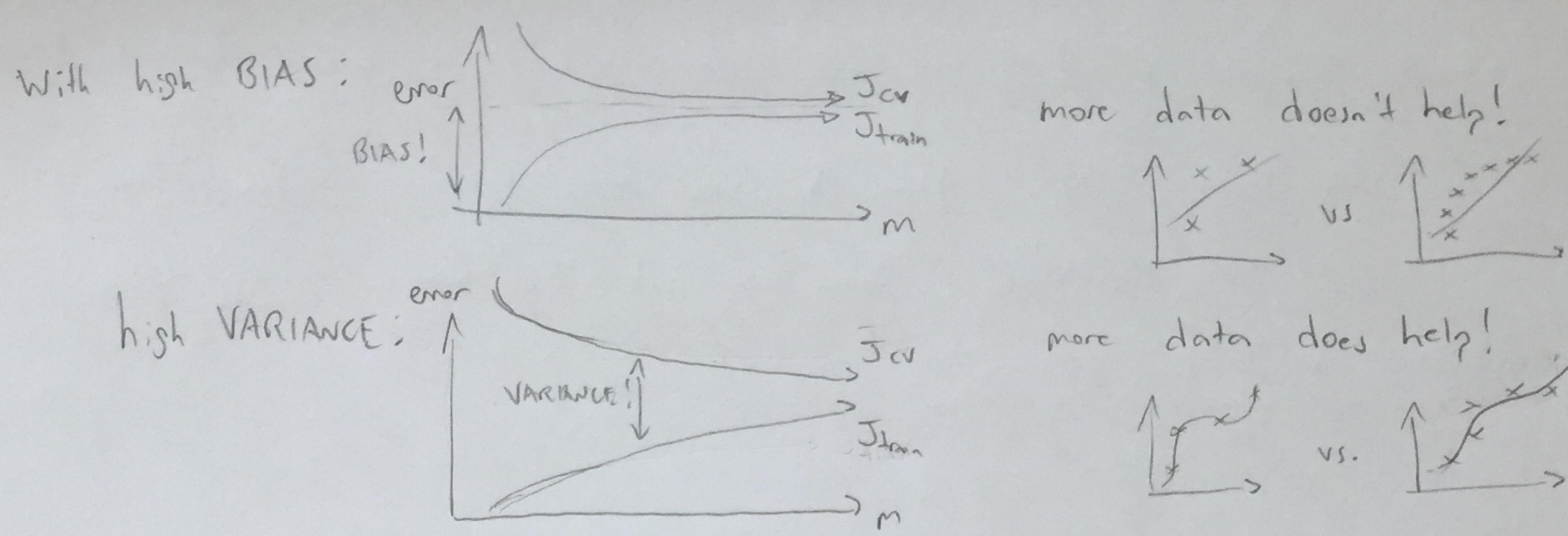


$$J_{train}(\theta) = \frac{1}{2m_{train}} \sum_{i=1}^{m_{train}} [h_\theta(x^{(i)}) - y^{(i)}]^2$$

plot

limit to m to get plot

$$J_{cv} = \text{" (w/ } M_{cv})$$



Summary:

- More training examples \rightarrow helps for high variance (look for $J_{cv} > J_{train}$)
- Smaller feature set \rightarrow fixes high variance (lower order d)
- Add features \rightarrow fixes bias (increase order)
- Changing $d \rightarrow \downarrow d$ for bias, $\uparrow d$ for variance fix
- For NN's: \Rightarrow low hidden units, low hidden layers \rightarrow "lower features"
 - \hookrightarrow prone to underfitting (bias) but computationally \downarrow
 - \Rightarrow more units / layers \rightarrow "higher features"
 - \hookrightarrow prone to overfitting (variance)
 - \hookrightarrow can use regularization to address overfitting (λ)!
 - \Rightarrow can compare NN's $J_{cv}(\theta)$!

Recommended Approach to Building Algorithms:

- ⇒ Start w/ simple algorithm → implement & test on CV. data
- ⇒ Plot learning curves (train+test error vs. # training examples / # features) to see what may help
- ⇒ Manually examine examples that algorithm makes error on → see if pattern
- ⇒ Words: → stemming software (N.L.P.) helps distinguish roots/meanings
 - ↳ can compare classification errors w/ or w/o stemming
 - upper/lowercase also compared
- ⇒ Always better to be quick n' dirty in initial implementation
 - ↳ helps decide where to spend time later
- ⇒ Skewed Classes:
 - ⇒ logistic / classification problem → find 1% test error
 - ↳ what if only 0.5% of patients have cancer?! → $y=0$ does better than your algorithm!
 - ↳ hard to tell if good model just by error...

Precision / Recall Metric:

		actual		⇒	Precision = $\frac{\# \text{ true pos}}{\# \text{ pred. pos} (\text{true + false +})}$	Recall = $\frac{\# \text{ true pos}}{\# \text{ actual pos} (\text{true + false -})}$
		1	0			
predicted	1	True +	False +			
	0	False -	True -			

- ⇒ Can't "cheat" by predicting only $y=0$, $y=1$

↳ algorithm w/ good precision & recall → good algorithm

Trading off Precision & Recall

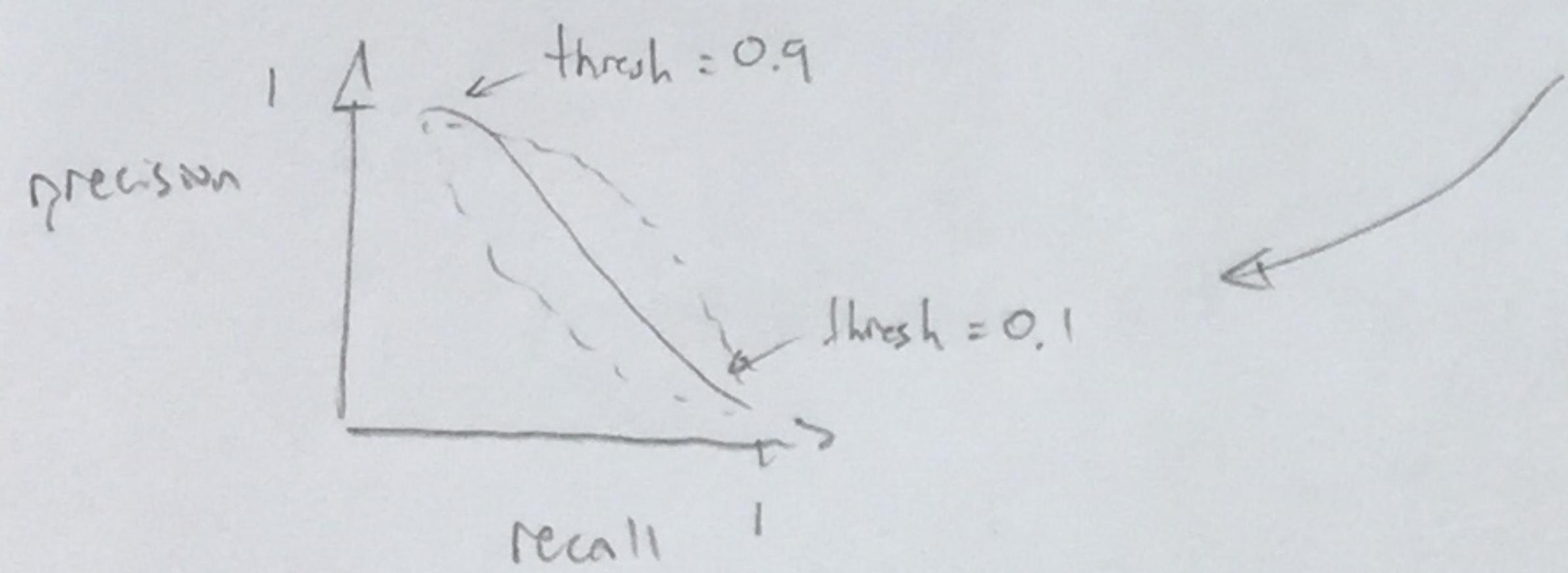
⇒ Say want to predict $y=1$ (e.g. cancerous) only if very confident

↳ change $h_0(x)$ log threshold from $0.5 \rightarrow 0.7$ (or 0.9)

⇒ HIGH PRECISION, LOW RECALL

⇒ vice versa, if want to be overcautious → $h_0(x)$ log threshold $0.5 \rightarrow 0.3$

⇒ HIGH RECALL, LOW PRECISION



⇒ How to compare algorithms given their prec./recall? (2 #'s)

	Precision:	recall:
Alg 1	0.5	0.4
Alg 2	0.7	0.1
Alg 3	0.02	1.0 ← predict $y=1$ all the time equivalent

↳ to get one comparison metric:

$$\boxed{F_1 \text{ Score:}} \quad F_1 = 2 \frac{PR}{P+R} \quad F_1 \in [0,1]$$

lower than either P or R

↳ compare F_1 scores of P & R on CV. set!

Large Data Rationale: (better to have more data assumption)

↳ many parameters → J_{train} small

↳ large training set → $J_{\text{train}} \approx J_{\text{cv}}$ ⇒ J_{cv} small!

⇒ features x must contain enough info to correlate/predict y , regardless of # parameters!