

ICCS482 Deep Learning

Lecture 5: underfitting, overfitting and regularization

Sunsern Cheamanunkul, Sep 22, 2020.

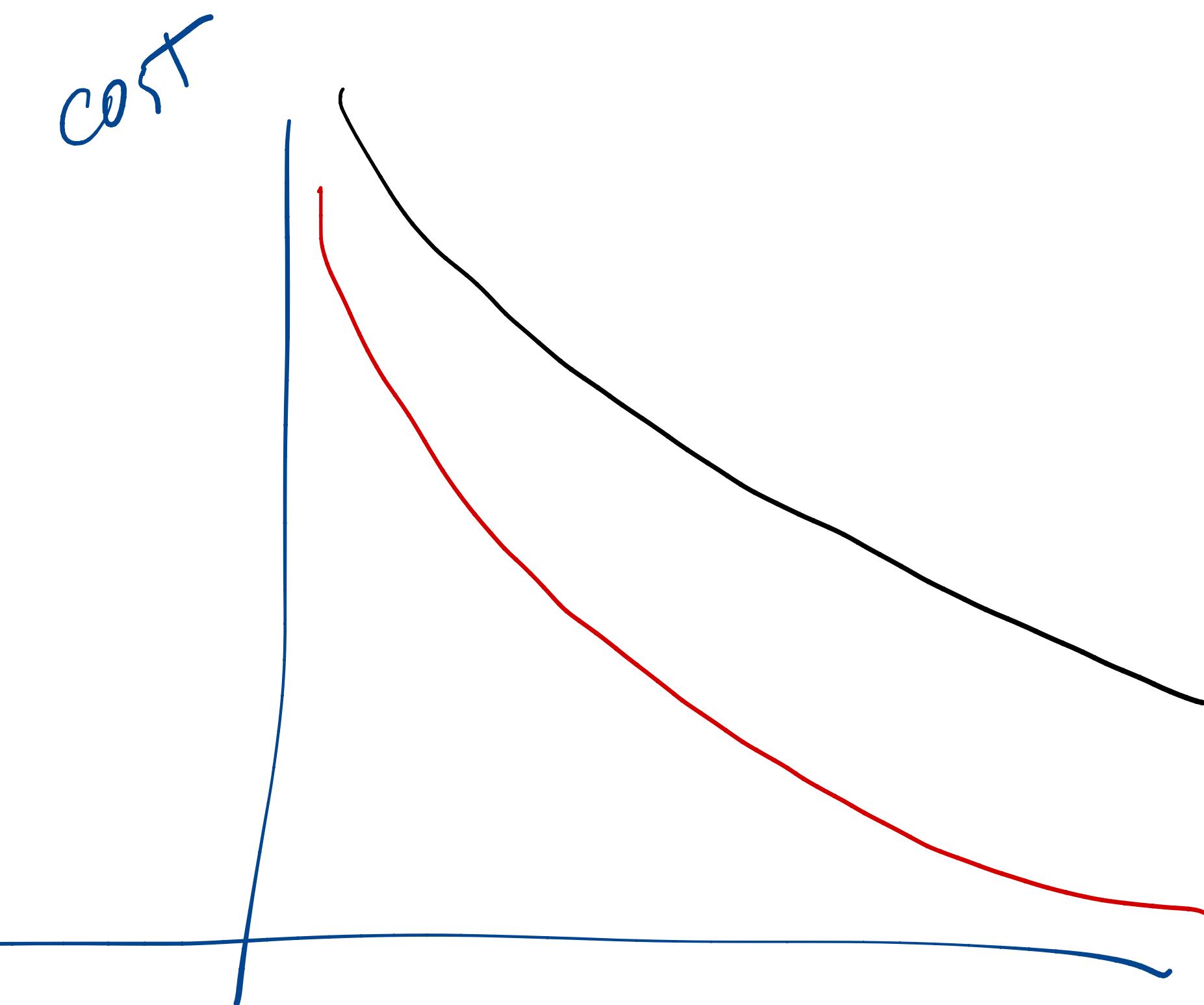
Diagnosing problems with NNs

Terms

- ① Training errors — Cost on Training Set
- ② Test errors — Cost on the test set

Scenario #1

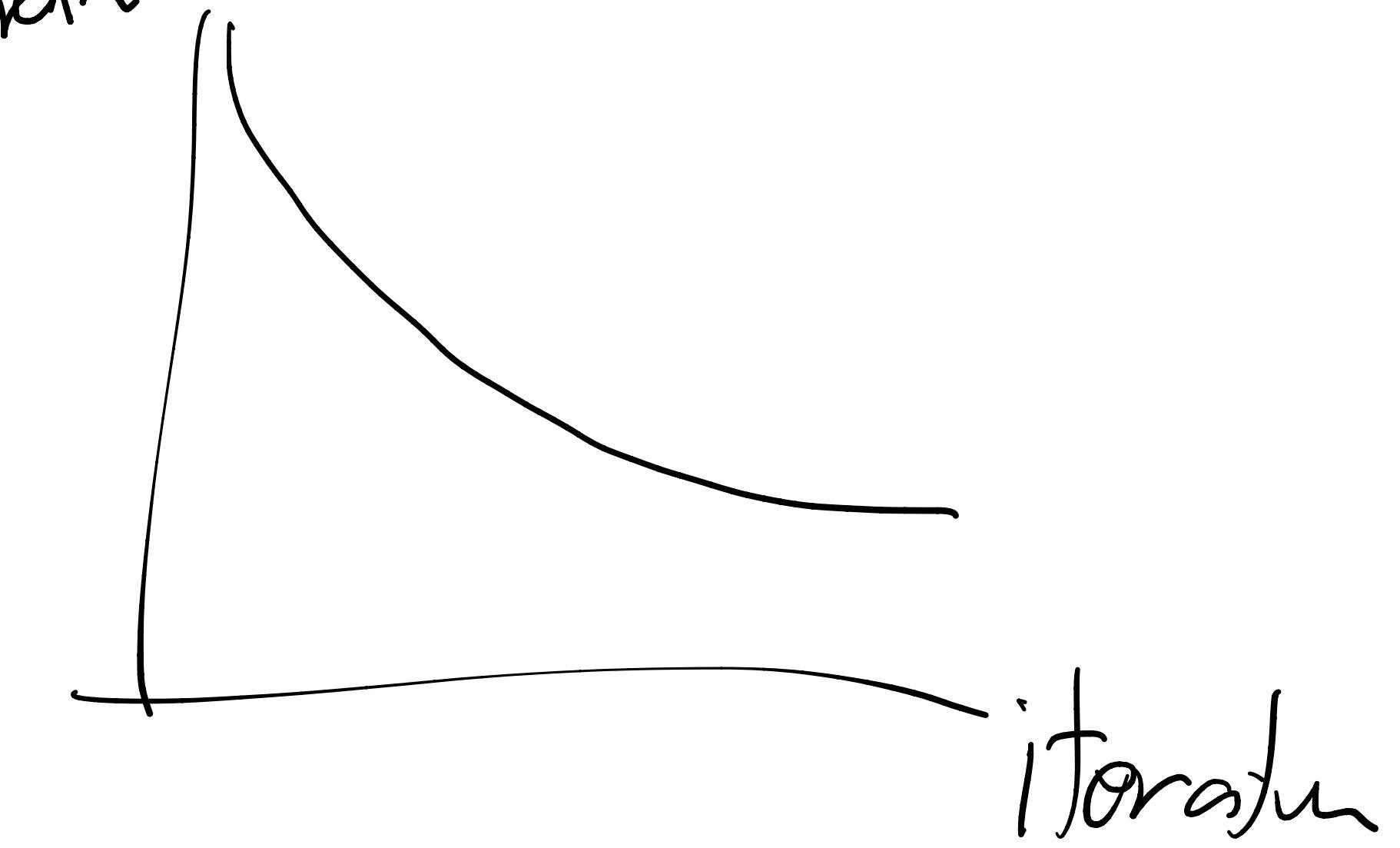
High training error



"underfitting"

(high cost on training set)

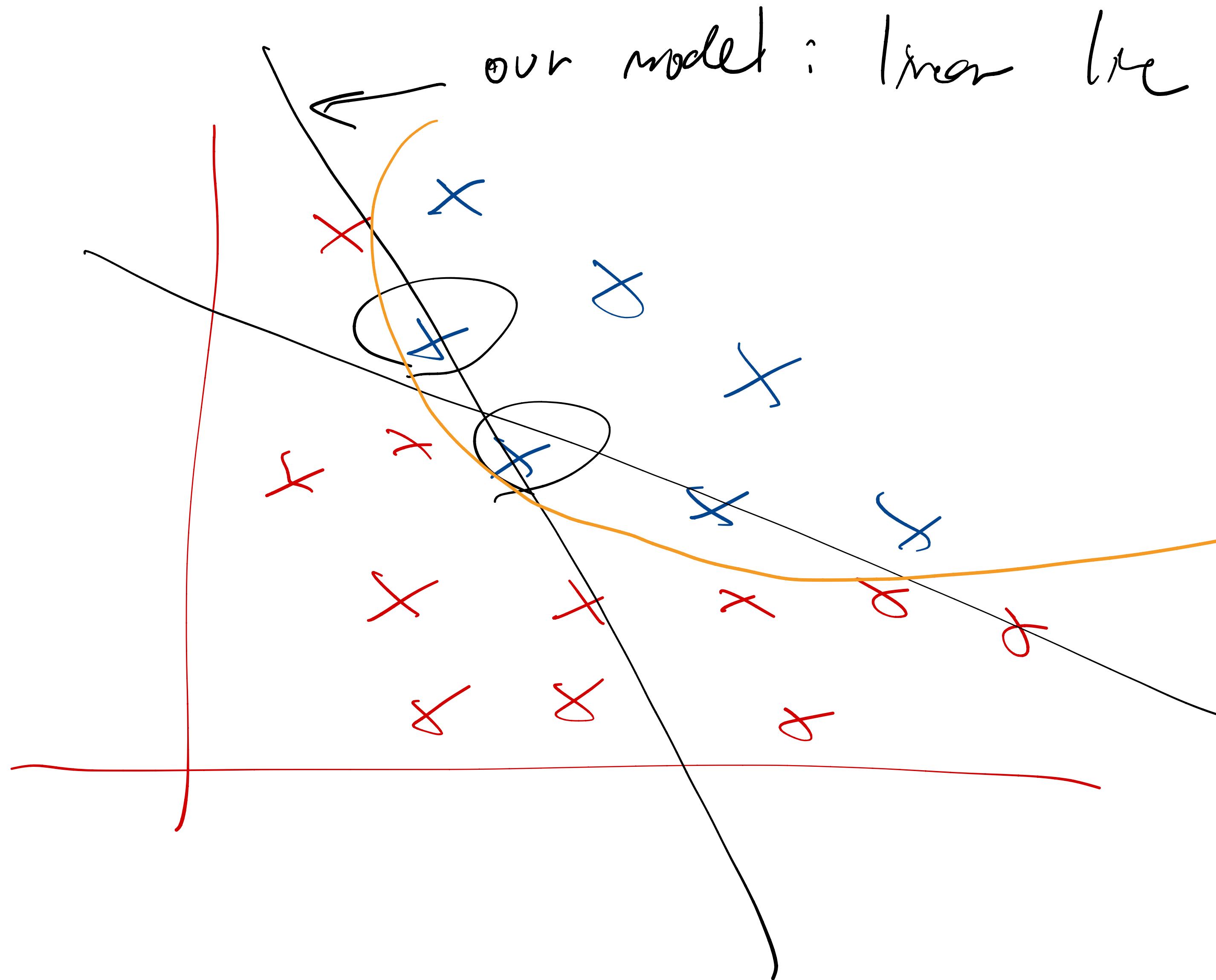
train error



iterations

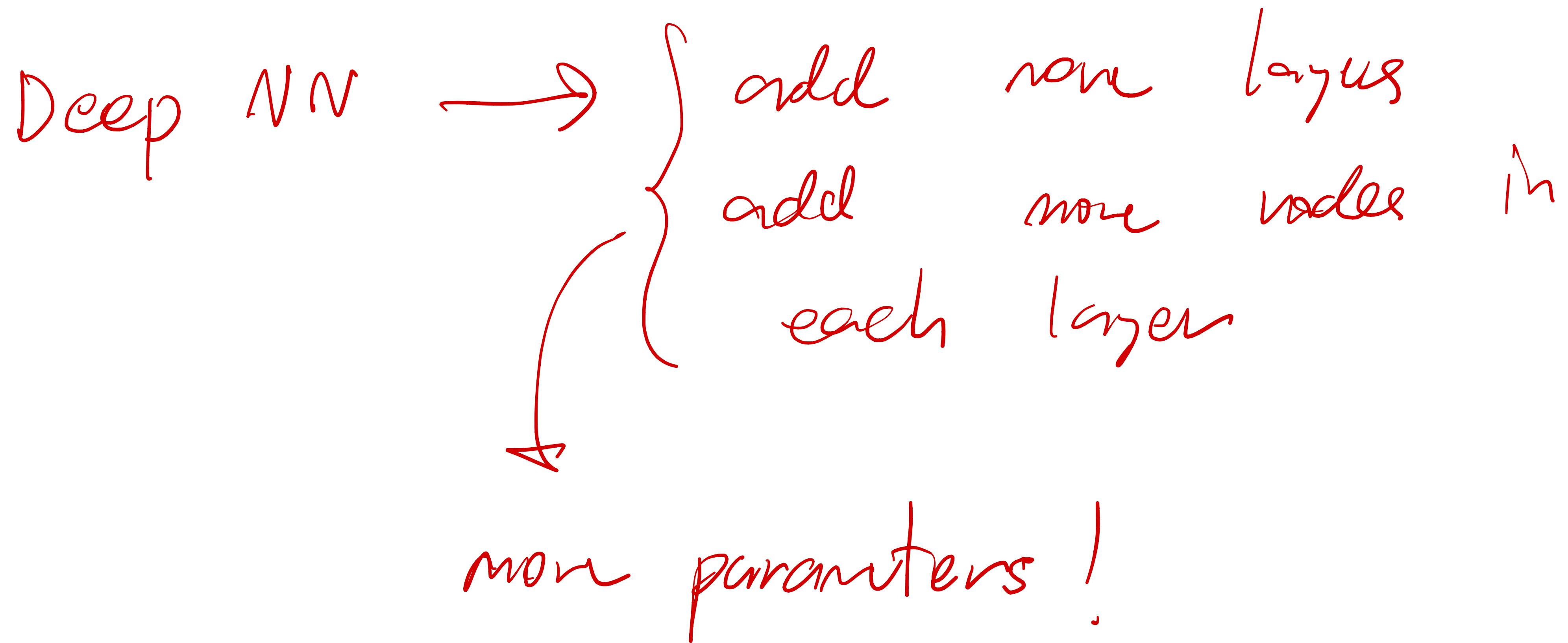
Why do we have high train error?

- ① model has not ~~enough~~ learned enough
- ② intrinsically very bad
(impossible to find a good role!)
- ③ model is not complex enough ~~or~~
(too simple)



↙ better model!
more complex
than a straight
line.

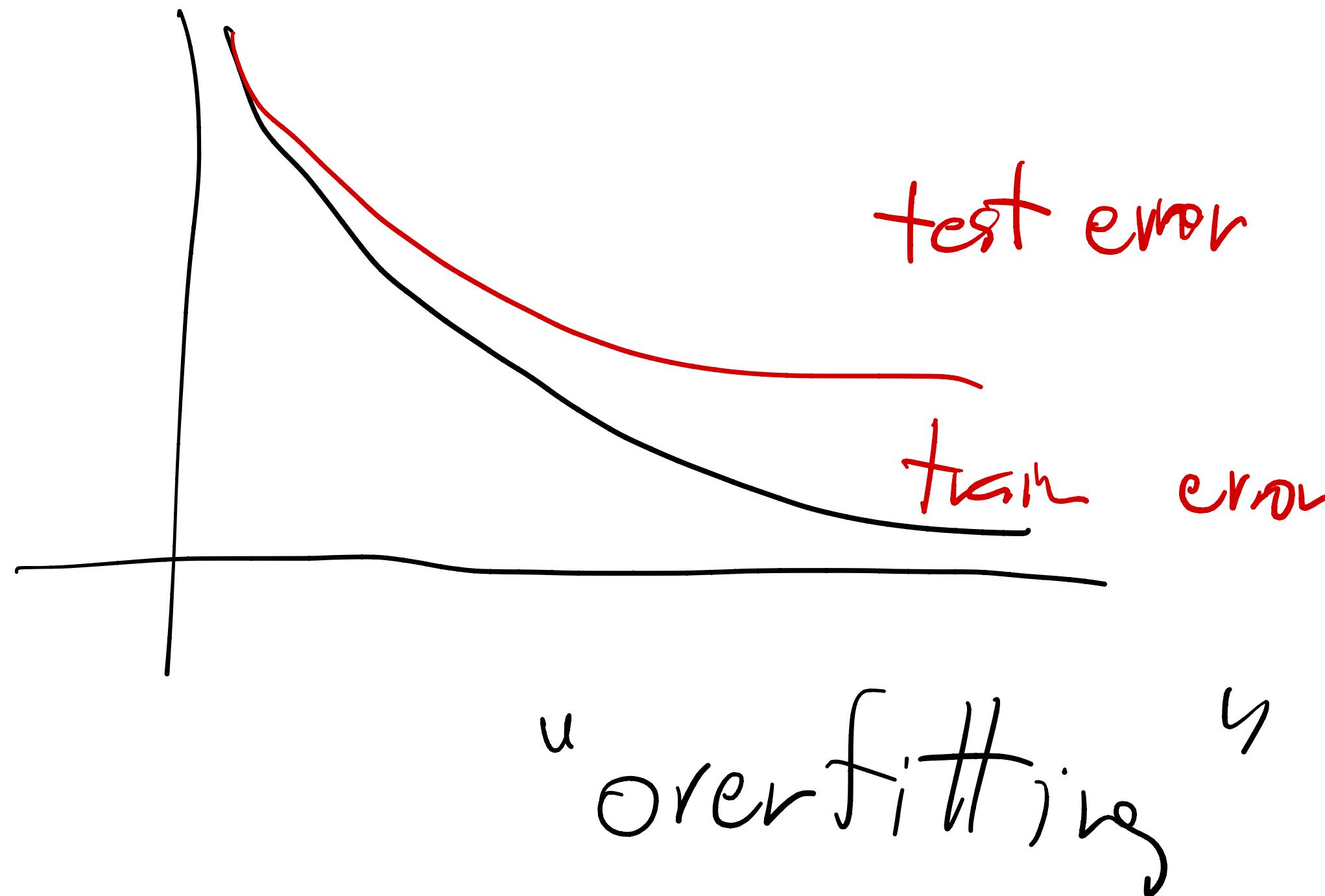
How to make a model more complex?



Scenario #2

Low train error

high test error



Reasons

- ① model too complex
- ② Train data and test data doesn't come from the same distribution
~ haven't seen enough.

Scenario #3

high train error] may be something wrong
low test error with our
train data

this is good.

(should not happen IRL)

scenario #4

high train err

high test error

model is not learning at all
- underfitting?

Underfitting

- ① try adding more layers
- ② try adding more hidden units
- ③ Add more parameters

Overfitting

- fit train data well
- doesn't work on test

the model generalises poorly!

What to do

- ① Decrease # of params
 - # of hidden layers / units
- ② Add more data - so the model sees more data
need variety!

③

Regularizations

- punishing the model in some way
- L₂ regularization:

$$J(\tilde{w}, \tilde{\zeta}) = \frac{1}{m} \sum_{i=1}^m l(y^{(i)}, \tilde{y}^{(i)}) + \frac{\lambda}{2m} \|\tilde{w}\|^2$$

Find $\tilde{w}, \tilde{\zeta}$ that minimizes $J(\tilde{w}, \tilde{\zeta})$ encouraged small weights.

$$\omega_1 = \begin{bmatrix} 100 \\ 200 \end{bmatrix}$$

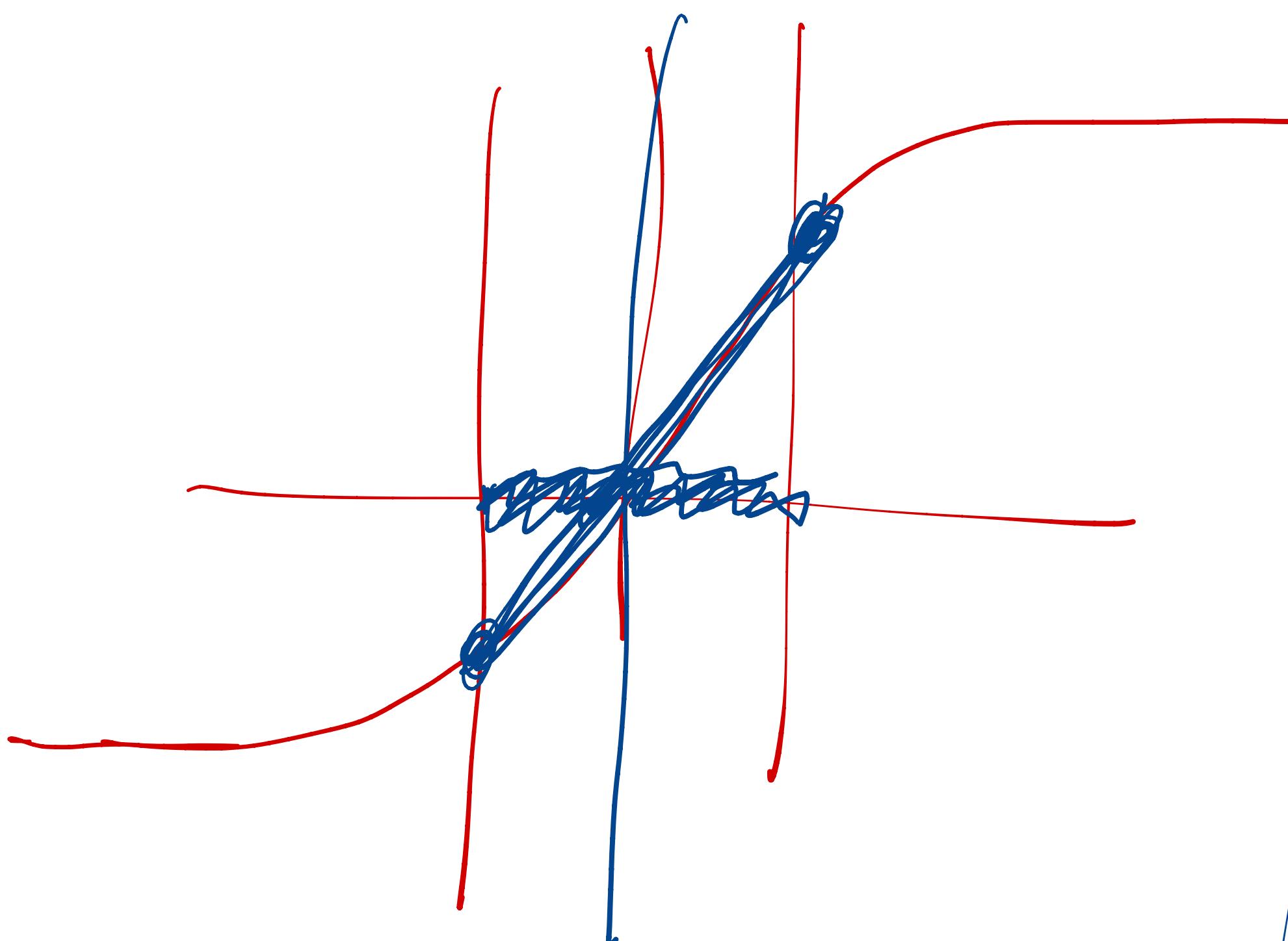
vs

$$\omega_2 = \begin{bmatrix} 0.1 \\ 0.2 \end{bmatrix}$$

this is problem by
 ℓ_2 regularizer

Why regularizers help with overfitting?

- ① make the model less complex



as if we have less

$$a^{(cl)} = g(z^{(cl)})$$

nodes
cls
w.o

z^(cl)

g(z^(cl))

as if we have less

Other regularizations

L_1 regularization

$$J(\tilde{w}, \tilde{s}) = \frac{1}{m} \sum_{i=1}^m l(s_i^{(i)}, \tilde{s}^{(i)}) + \lambda \|\tilde{w}\|,$$

we want L_1 of w

$$\|\begin{bmatrix} 3 \\ -3 \end{bmatrix}\|_1 = |3| + |-3| = 6$$

to be small

\rightarrow we want

$$\|\begin{bmatrix} 3 \\ 1 \end{bmatrix}\|_2 = \sqrt{3^2 + 1^2} = \sqrt{10}$$

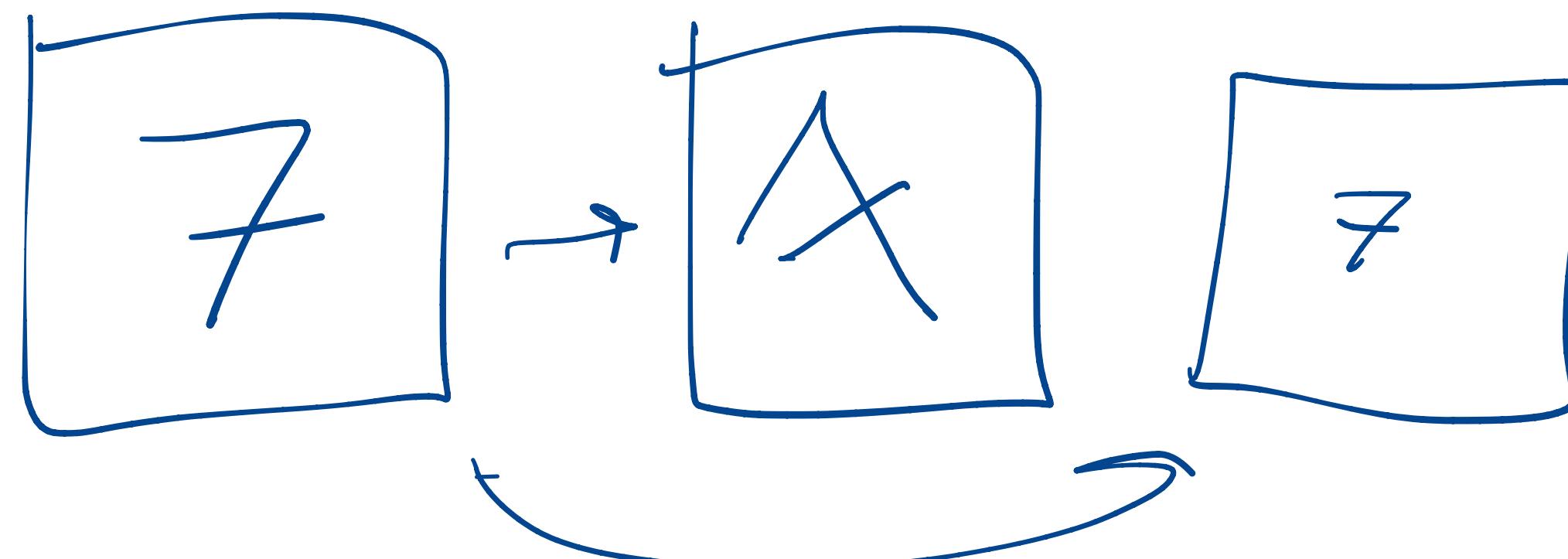
0 to appear often in w

Dropout regularization

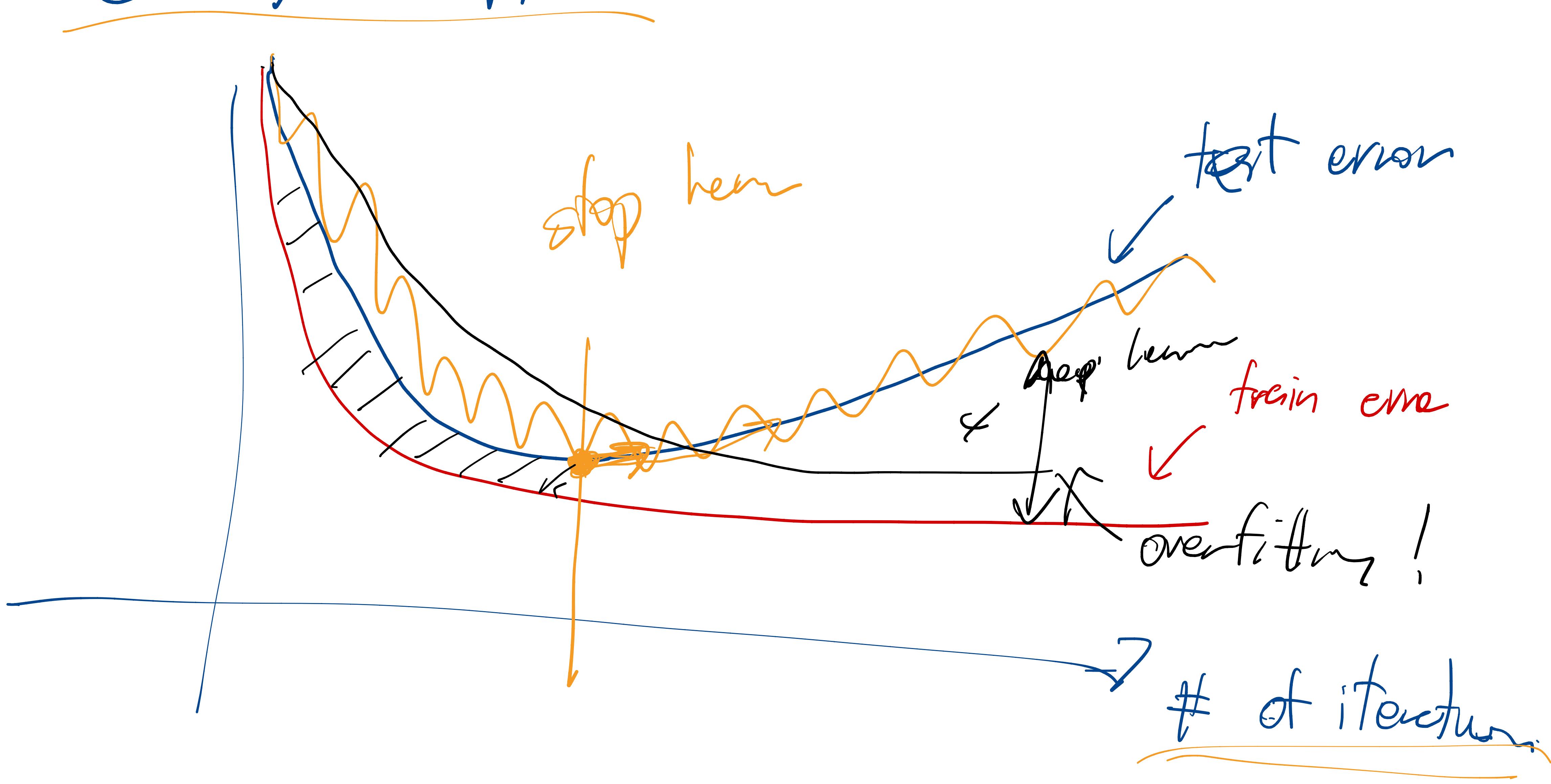
- randomly zero out nodes
- do not want to frost any single nodes too much
- hope this helps with generalization
- don't drop during test phase

Other methods to fight overfitting

- Add more data!
 - sometimes this is not possible.
 - data augmentation techniques
 - take existing examples and
freak them a little e.g.



• Early Stopping



6