

Last time:
$$MSE_{\text{train}} = \frac{1}{N} \sum_{n=1}^N (y_n - \hat{y}_n)^2$$

training data: $\{(y_n, x_n)\}_{n=1}^N$

$\hat{y}_n = \hat{f}(x_n)$

fit \hat{f}

What do I care about when evaluating a method.

Things I don't care about:

performance of \hat{f} on the training data

↳ I already know the answer: y_n

Things I care about:

performance of \hat{f} on new data I haven't seen before.

↳ I don't know the answer

Defn:

training data: data fed to approach to create \hat{f} ← function of training data

$$\{(x_n, y_n)\}_{n=1}^N$$

test data: data from the same data generating process but that wasn't used to "train" \hat{f} .

$$\{(x_n^*, y_n^*)\}_{n=1}^{N^*}$$

n

Scheme for evaluation of \hat{f}

- ① use training data to train \hat{f}
- ② evaluate \hat{f} on the test data

double up on metrics

RSS_{train}

RSS_{test}

MSE_{train}

MSE_{test}

$RMSE_{\text{train}}$

$RMSE_{\text{test}}$

calc. on training

replace training data w/
test

Why? Evaluating on training data is potentially misleading. (typically)

Intuition: Circular argument.

$$(*) \quad \hat{f} = \underset{f}{\operatorname{argmin}} \mathbb{E}[L(f)]$$

\hat{f} is chosen to minimize the expected (empirical) risk.

I train \hat{f} , by finding a \hat{f} that predicts my

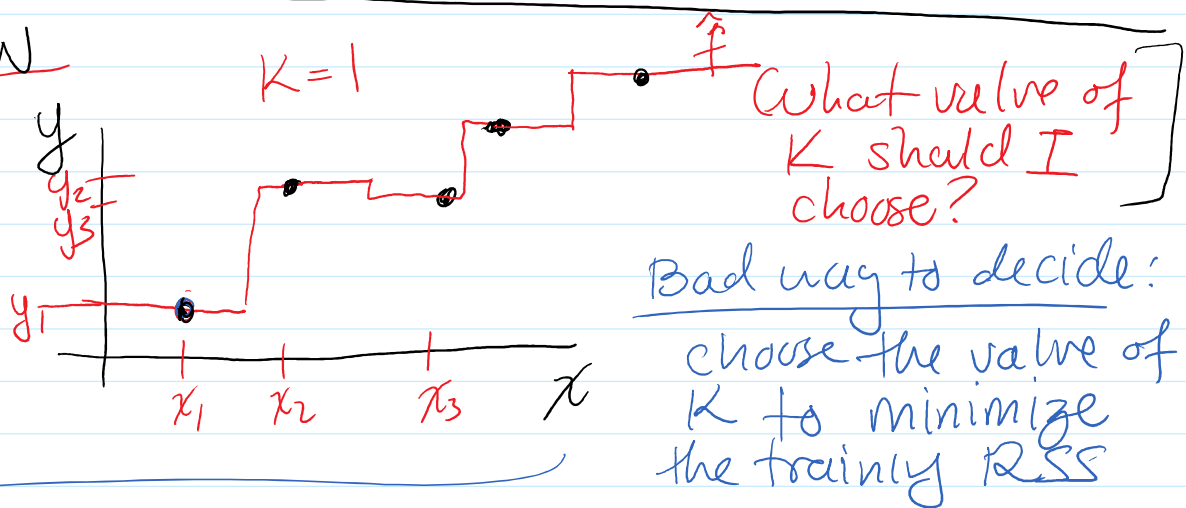
training 'data' well. ~

So.... \hat{f} predicts training data well.

Result: training metrics are too optimistic on performance of \hat{f} on new/independent data

i.e. $RSS_{\text{train}} < RSS_{\text{test}}$

Ex. KNN



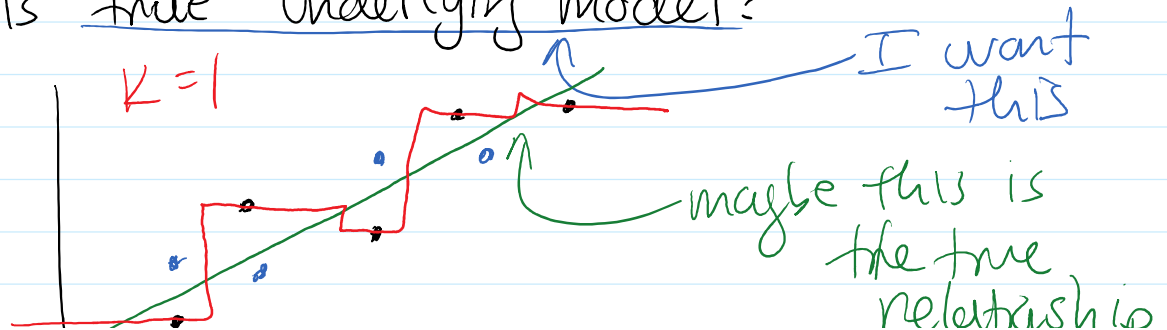
Q: what value of K minimizes RSS_{train} ?

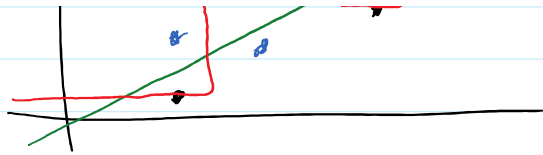
$K=1$

$$RSS_{\text{train}} = \sum_{n=1}^N (\hat{y}_n - y_n)^2 = 0$$

b/c $\hat{y}_n = y_n$

What is true underlying model?



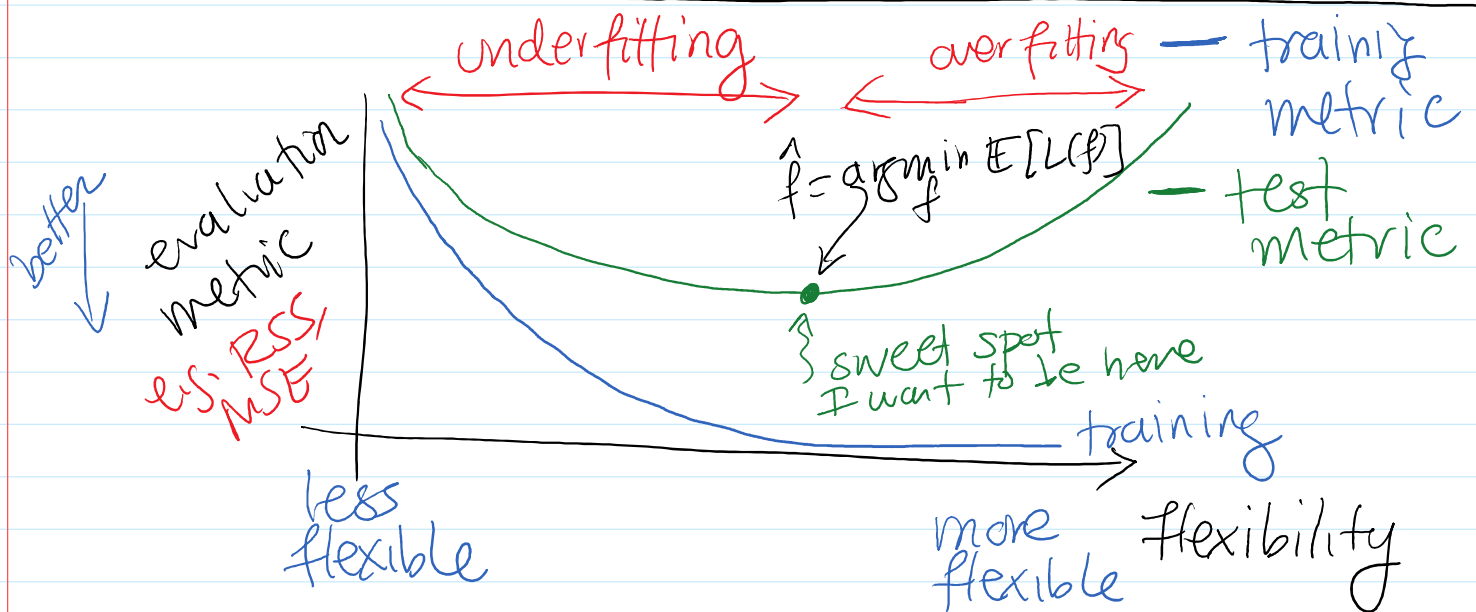


the true relationship

I'm not following the true underlying trend,
I'm chasing the noise.

when I generate new data:
it won't follow noise, but the underlying trend.

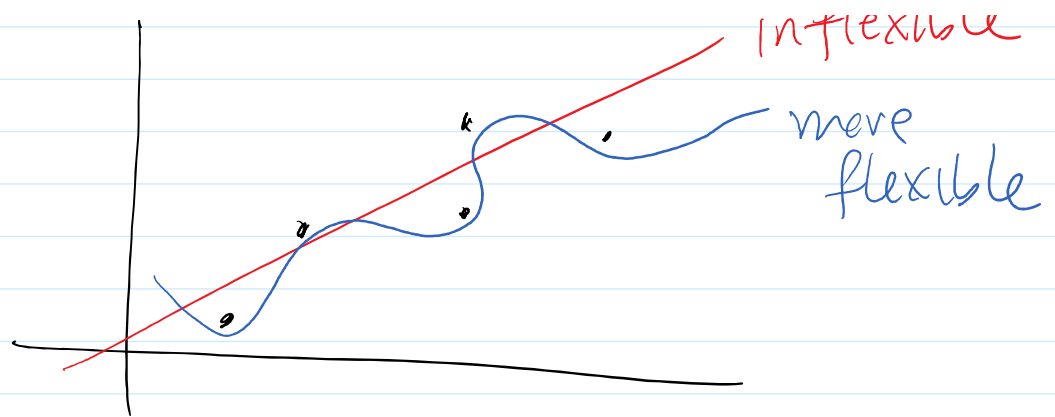
We call this overfitting.



- KNN w/ high $K = N$
- lin. reg. w/ low p
(not many covs.)

- KNN w/ low $K = 1$
- lin. reg. w/ lots of added covs.
high p

inflexible



How do we solve this in real life?

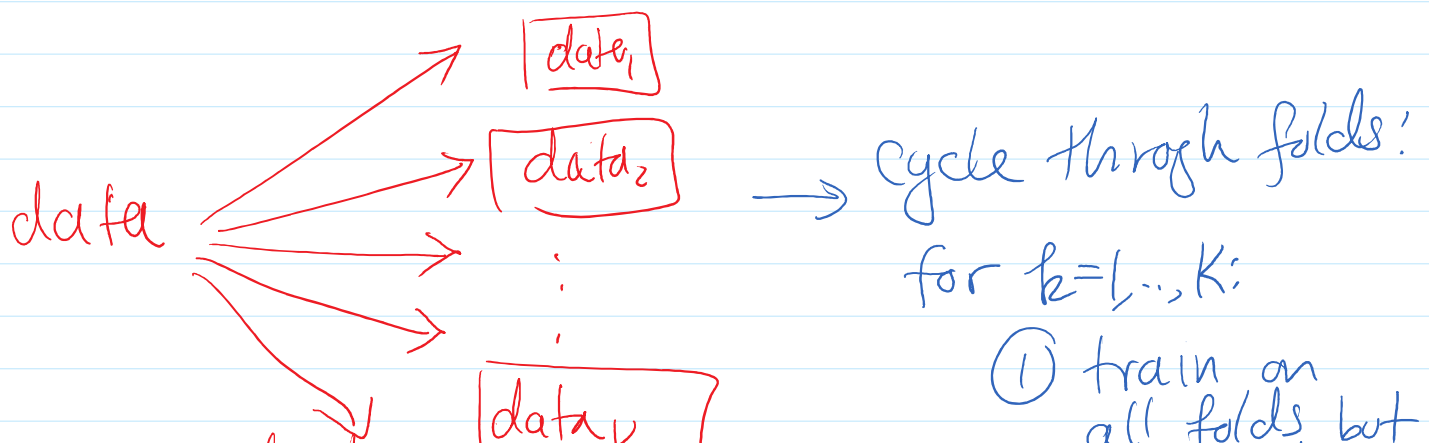
I get some data set:

→ want to use to make f to predict well

Validation Set: test/train split

data $\xrightarrow{p\%}$ training data (fit)
 $\xrightarrow{(1-p)\%}$ test data (evaluate)
 typically $p \approx 50\%$
 split (randomly)

Cross-Validation: K-fold cross-validation



randomly
split in "folds"

data_k

(1) train on
all folds but
 k^{th}

(2) evaluate on
 k^{th} fold

→ after this I get K metrics

m_1, m_2, \dots, m_K

summarize into a final metric

e.g. $m = \text{median}(m_1, \dots, m_K)$.
