# Lecture 5: Evaluation

How do I choose among the various models?

We need some way to evaluate which model building process is going to perform better.

→ use OLS w/ 5 vars
→ use kNN w/ $k = 5$
→ use kNN and optimize over possible values of $k$

A (maybe not always great) way of doing this!

Calculate some performance metric on $\hat{f}$ from the training data

e.g.

① training residual sums of squares

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

② training mean sq. error

$$MSE_{train} = \frac{RSS_{train}}{N}$$

③ training root mean sq. err: $RMSE_{train} = \sqrt{MSE_{train}}$

(4) training $R^2_{train} = 1 - \dfrac{RSS_{train}}{TSS_{train}}$

$TSS_{train} = \sum_n (y_n - \bar{y})^2$

$\uparrow$ % of var. explained by $\hat{f}$

Why isn't a training metric always a good measure of performance?

$\rightarrow$ I don't actually care about performance on my training data
(I already know the answer)

$\rightarrow$ I actually care about is performance of $\hat{f}$ on new/unseen data.
(generalization performance)

ERM:

$$\frac{1}{N} \sum_{n=1}^{N} L(y_n, f(x_n)) \approx \mathbb{E}\left[ L(Y, f(X)) \right]$$

avg. loss over training data

actual loss

Similarly:

$$Metric(\{y_n, f(x_n)\}) \approx \mathbb{E}\left[ Metric(\{Y, f(X)\}) \right]$$

training metric/perf.

generalization performance

Focusing too much on training metric can be misleading.

This happens b/c when I evaluate $\hat{f}$ on my training data I am evaluating it on the same data used to build $\hat{f}$.

So my model building process has already seen the traing data — its not a fair evaluation.
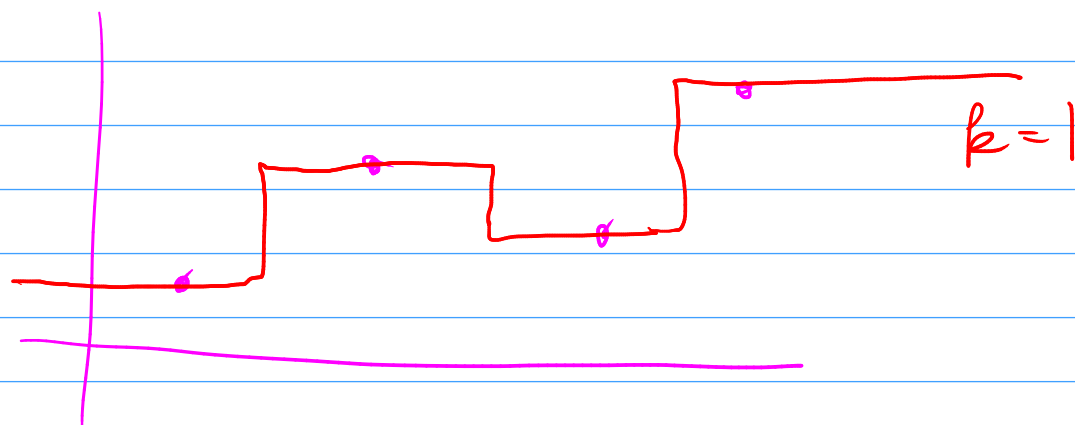
Ex. This can happen in SML too.

Consider optimizing the value of $k$ for kNN by choosing the value of $k$ that minimizes $RSS_{train}$.

What value of $k$ do I choose?
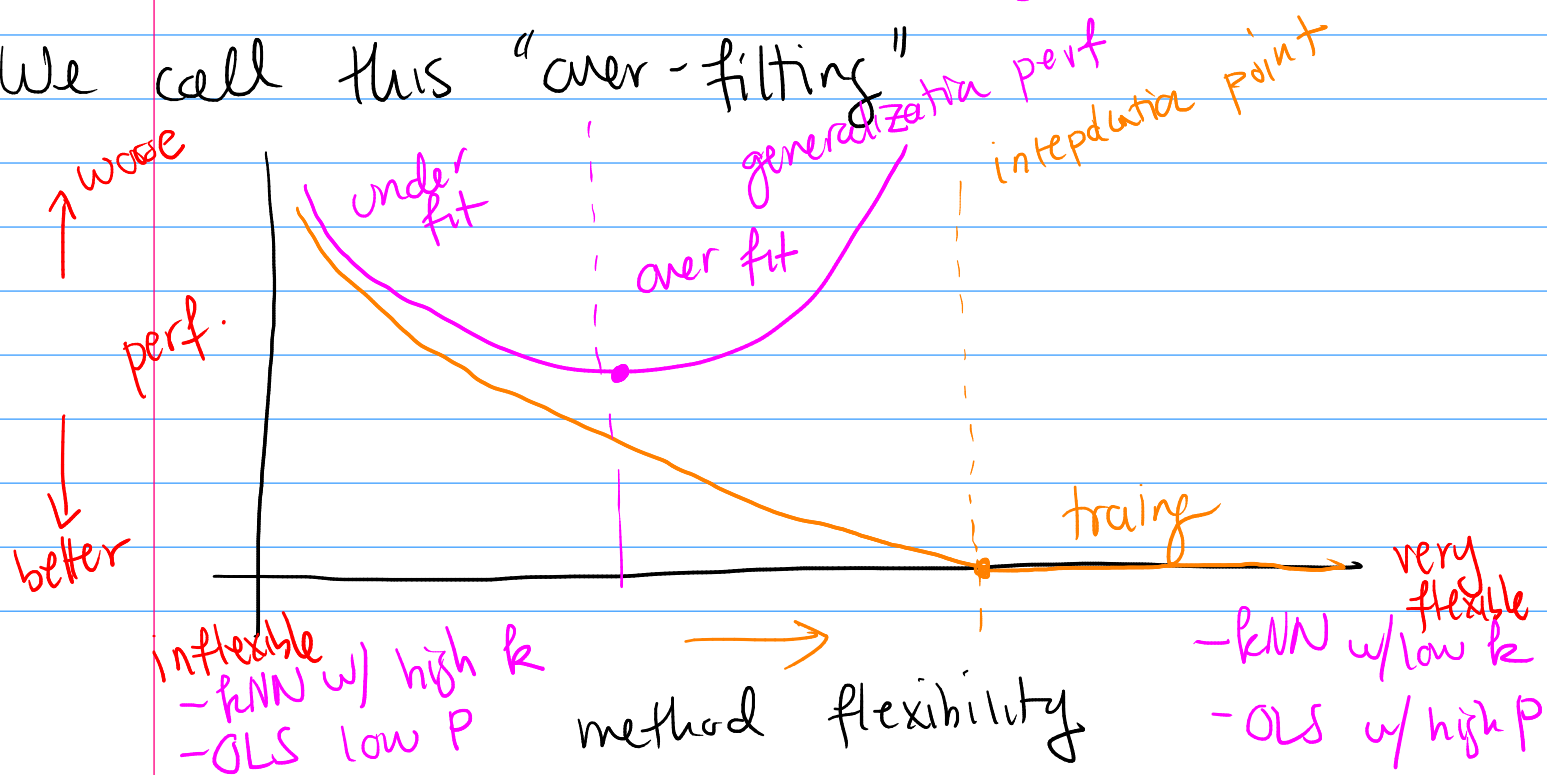
$$k = 1$$

When $k=1$ I interpolate my traing data



$$k=1$$

So $\hat{y}_n = y_n$ for traing data $n=1,\ldots,N$

hence $RSS_{train} = \sum_n \underbrace{(y_n - \hat{y}_n)^2}_{0} = 0$

- - - - - - - - - - - - - - - - - - -

A similar story is true for OLS when choosing what terms to include in the design.

as $p \to N$ I will interpolate my traing data

We call this "over-fitting"



worse

perf.

better

generalization perf

interpolation point

under fit

over fit

traing

very flexible

inflexible
— kNN w/ high $k$
— OLS low $P$

method flexibility

— kNN w/ low $k$
— OLS w/ high $P$

How do we solve this?

Let's estimate generalization perf.

Need independent set of (unseen) data to
evaluate our model on.

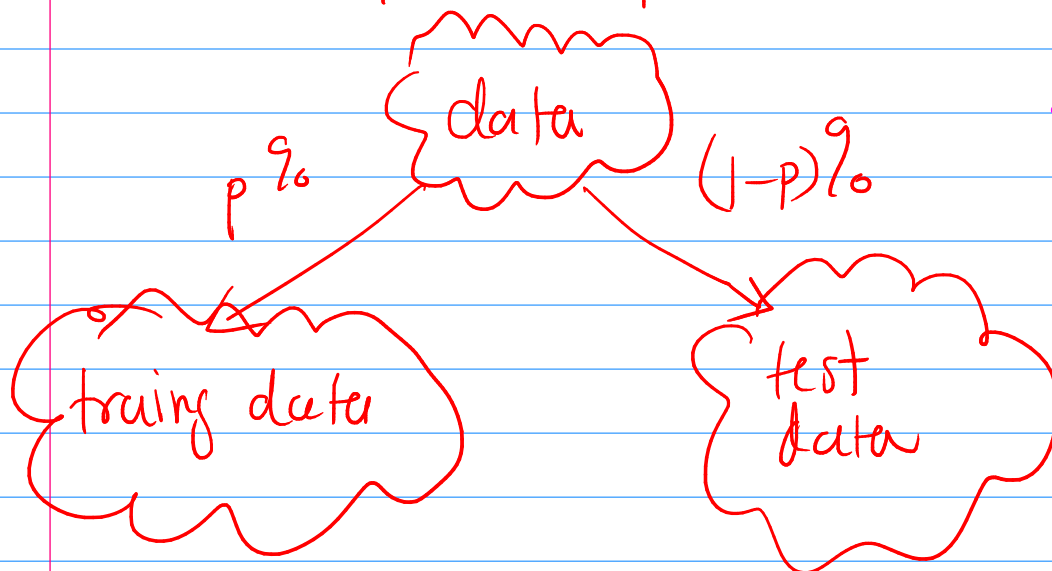Testing Data: $\{(\underset{\sim}{x}_{test,n}, y_{test,n})\}_{n=1}^{N_{test}}$

Procedure: ① use traing data to build $\hat{f}$

② eval. perf. of $\hat{f}$ on testing data

$RSS_{test}, MSE_{test},$
$RMSE_{tot}, R^2_{test}$

↳ an estimate of generalizatia
perf. of my model buildij
process.

How do I get testing data?

Do a test/train split.



$p$ %         data         $(1-p)$%         $p = 90, 80, 50$

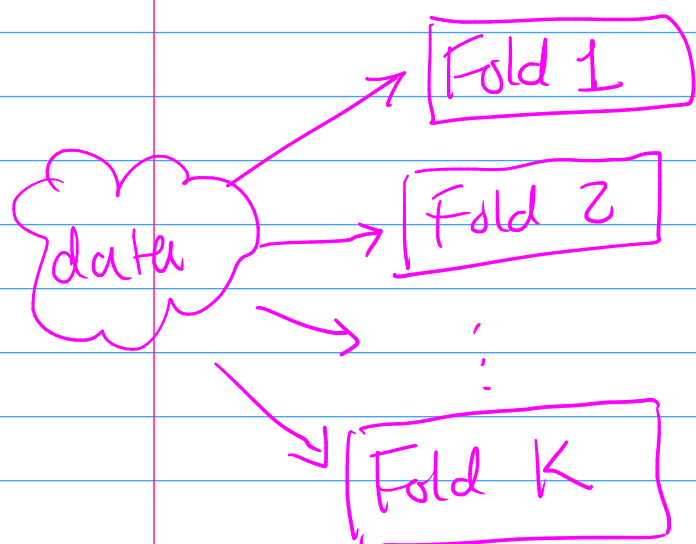traing data                    test
data

One can do this split multiple times which is called cross-validation

## K-Fold Cross Validation



Cycle through folds

For $k = 1, ..., K$

① train model on all but fold $k$

② test model on fold $k$
call this test metric $m_k$

So at the end of loop I have

$$m_1, m_2, ...., m_K$$

can combine to get overall perf. metric

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

But wait! If I train $K$ different models which do I use?

None. I train a model using all of my data.

We want to evaluate perfomance generally
for two reasons:

(1) I want to know how good this
MBP might be

(2) I want to choose among varias
models.

We can use a test/train split or X-validation
to choose among models — but we need to
be very careful.