

Lecture 5: KNN regression and Evaluation

KNN regression:

$$\hat{f}(\underline{x}) = \frac{1}{K} \sum_{n \in N_K(\underline{x})} y_n = \text{avg. training } y_n \text{ s whose corresp. } \underline{x}_n \text{ s are one of the } K \text{ nearest neighbors to } \underline{x}$$

What happens as I change K ?

General rule: K control model flexibility

Small $K \rightarrow$ very flexible methods

Large $K \rightarrow$ very inflexible methods

how complicated
my predictor
fn is

Evaluation

How do I choose among my various models?

e.g. choose K for KNN

e.g. choose design for OLS

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

Naive way: Calculate some performance metric of my model \hat{f} on the training data

e.g. (1) training RSS

$$RSS_{\text{train}} = \sum_n (y_n - \hat{y}_n)^2$$

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

(2) training MSE

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

(3) training RMSE

$$RMSE_{\text{train}} = \sqrt{MSE_{\text{train}}}$$

(4) training R^2

$$R^2_{\text{train}} = 1 - \frac{RSS_{\text{train}}}{TSS_{\text{train}}}$$

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

\approx % of var. explained by \hat{f}

Why shouldn't we do this?

(Why isn't a metric calc. on training good?)

→ I don't actually care about predictions on my training data.

→ I actually care about performance on new/unseen data.

[Generalization Performance]

Focusing too much on training data is misleading i.e. - a bad surrogate for generalization performance.

Ex. If I want to choose a "best" value for K in KNN what happens if I choose the val. to minimize RSS_{train} ?

Always going to choose $K=1$ b/c then $\hat{y}_n = y_n$ (Interpolate) so $RSS_{\text{train}} = 0$.

Similarly for OLS if I let

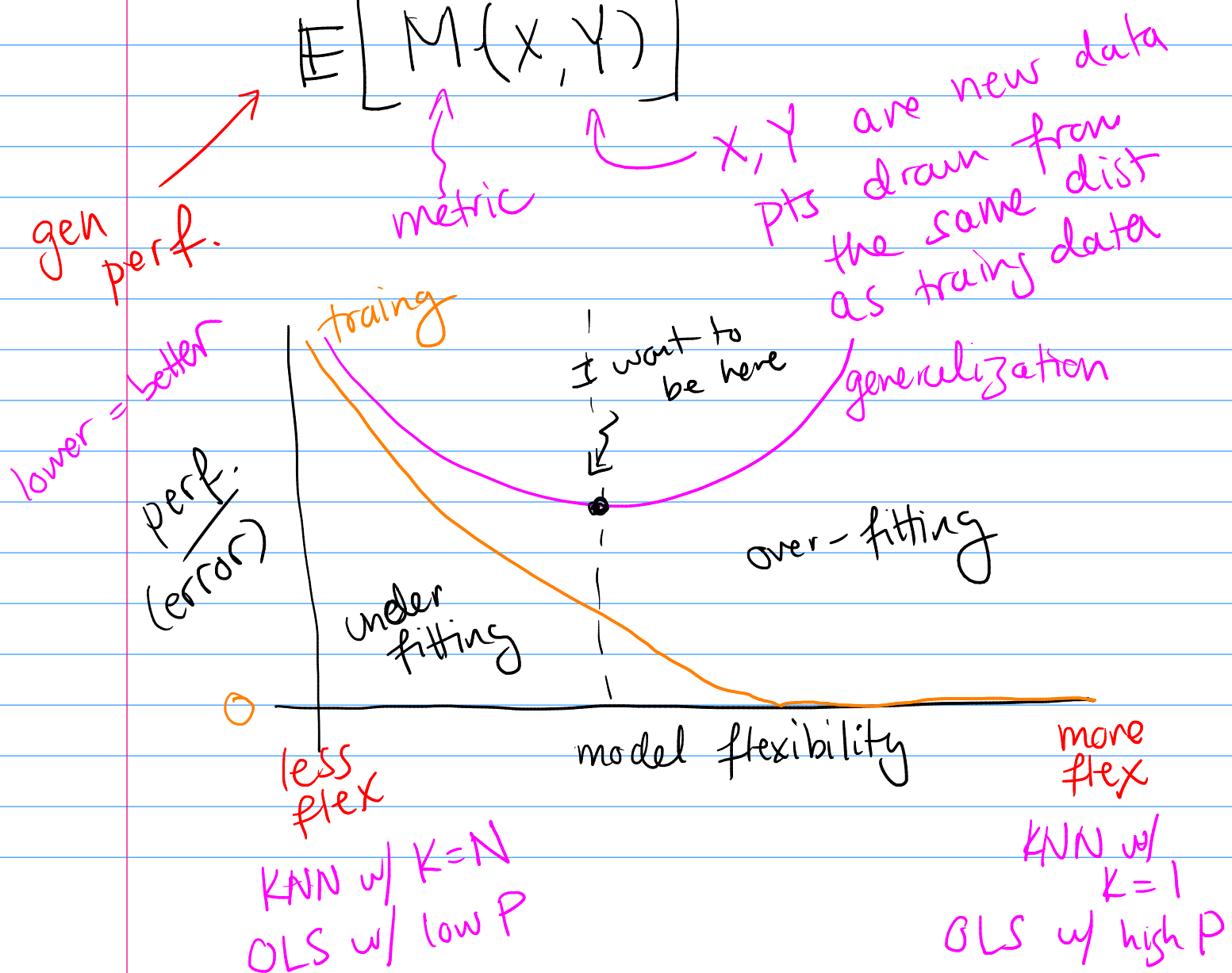
$$P \rightarrow N$$

I will eventually Interpolate my training data.

[If I fit a $N-1$ th degree poly. to N pts, I will interpolate]

What we really want to calculate is

$$\mathbb{E}[M(X, Y)]$$



How do we do this?

I need some data my approach (MBP) hasn't seen before

We'll call this test data: $\{(x_{n,\text{test}}, y_{n,\text{test}})\}_{n=1, \dots, N_{\text{test}}}$

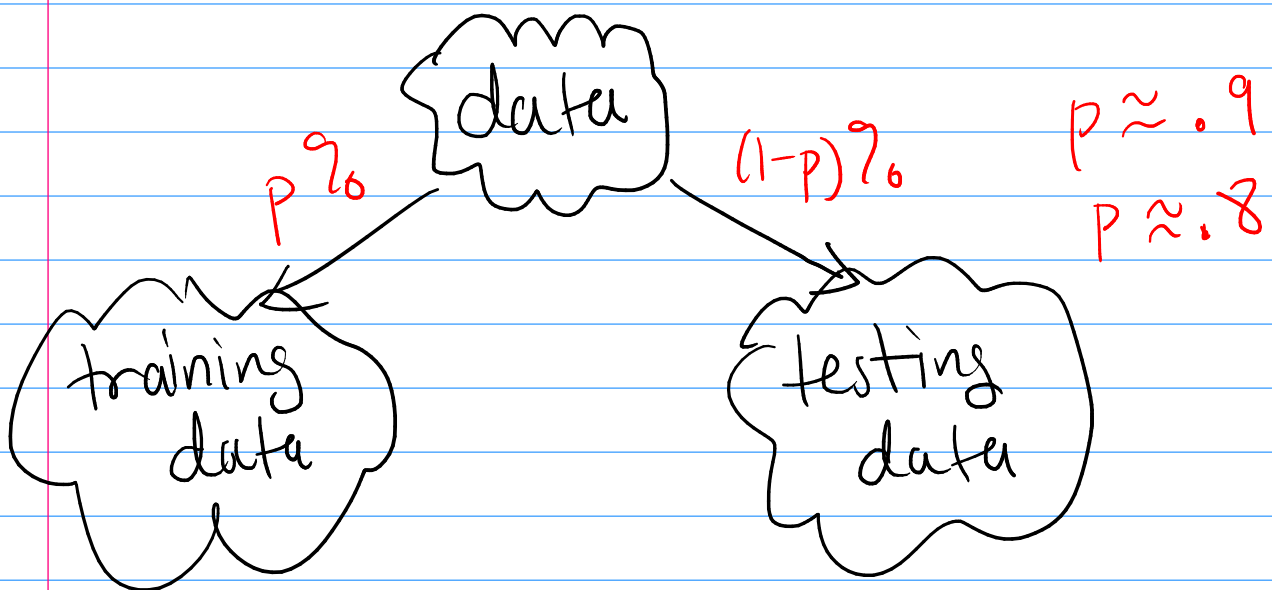
$$E[M(x, y)] \approx M(\text{test data})$$

Procedure:

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

How do I get test data?

Do a train/test split

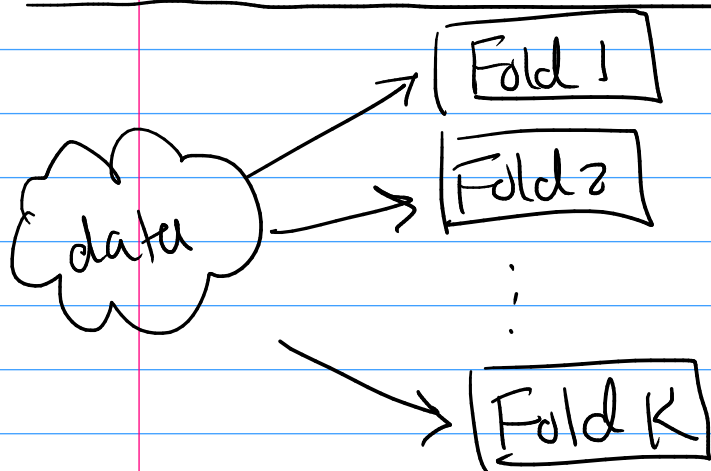


Can also do this multiple times,

Cross-Validation

K-Fold X-Validation

e.g. $K=10$



For $k=1, \dots, K$

① train model on all but k^{th} fold

② evaluate on the k^{th} fold

let m_k = perf. metric for this iter

At the end I have

$m_1, m_2, m_3, \dots, m_K$

Can combine together to get an overall metric e.g.

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

OK, so now I have K different models, which do I use?

None. I should probably train using all my data.