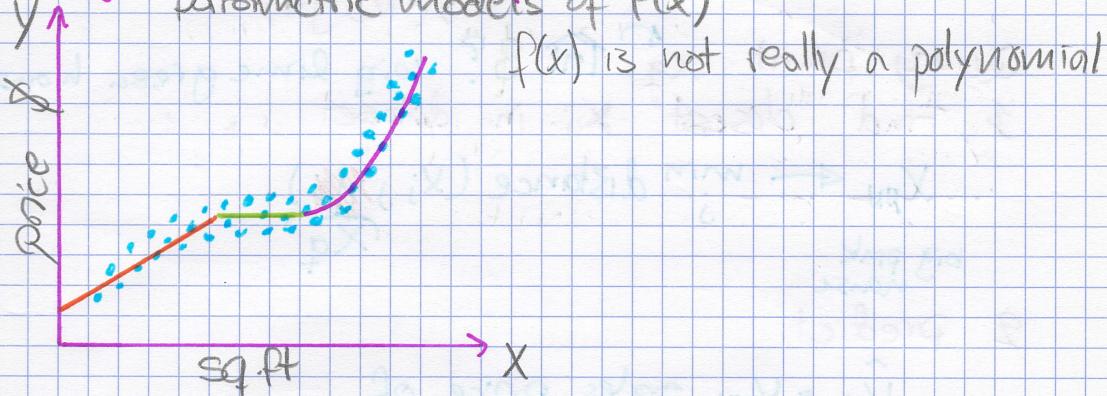


Nearest Neighbors & Kernel Regression

Limitations of parametric regression

Fit globally vs. Fit locally

parametric models of $f(x)$



What alternative do we have?

If we:

- Want to allow flexibility in $f(x)$ having local structure

- Don't want to infer "structural breaks"

What's a simple option we have? . . .

- Assuming we have plenty of data ...

1 - Nearest neighbor regression approach

Fit locally to each data point

Predicted value = "closest" y_i



1 nearest neighbor
(1-NN)
regression

1-NN regression more formally

Dataset of (Int , $\$$) pairs: $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$

Query point: x_q Int $\$$? $\text{big lime green house}$

1. Find "closest" x_i in dataset

$$x_{NN} \leftarrow \min_j \text{distance}(x_j, x_q)$$

big pink house

2. predict

$$\hat{y}_q = y_{NN} \text{ sales price of } \cancel{x_{NN}} \text{ big pink house}$$

Visualizing 1-NN in multiple dimensions



Voronoi tessellation
(or diagram)

- Divide space into N regions, each containing 1 data point

- Defined such that any x in region is "closest" to region's data point

x closer to x_i
than any other x_j
for $j \neq i$

Distance metrics:
Defining notion of "closest"

In 1D, just Euclidean distance:

$$\text{distance}(x_j, x_q) = |x_j - x_q|$$

In multiple dimensions,

- can define many interesting distance functions
- most straightforwardly might want to weight different dimensions differently

Weighting housing inputs
Some inputs are more relevant than others

Scaled Euclidean distance
Formally, this is achieved via

$$\text{distance}(x_j, x_q) =$$

$$\sqrt{a_1(x_j[1] - x_q[1])^2 + \dots + a_d(x_j[d] - x_q[d])^2}$$

↑ weight on each input (defining relative importance)

bed rooms

bathrooms

sq ft. living

sq ft. lot

floor s.

year built

year renovated

waterfront

Other example distance metrics

Mahalanobis, rank-based, correlation-based, cosine similarity, Manhattan, Hamming, ...

1-NN algorithm

Performing 1-NN search

- Query house:

- Dataset:



- Specify: Distance metric

- Output: Most similar house

Initialize $\text{Dist2NN} = \infty$, = \emptyset

For $i = 1, 2, \dots, N$

* closest house

Compute $\delta = \text{distance}(\text{house } i, \text{query house})$

if $\delta < \text{Dist2NN}$

set =

set $\text{Dist2NN} = \delta$

Return most similar house

* closest house

to query house

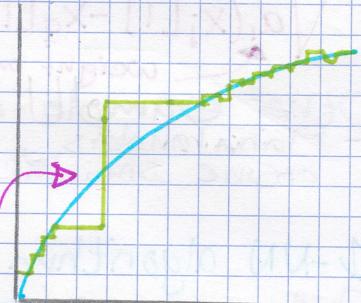
Nearest Neighbors Kernel $k=1$



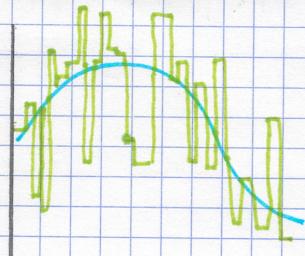
Fit looks good for data dense in x and low noise

Sensitive to regions with little data

Not great at interpolating over large regions



Also sensitive to noise in data



Fits can look quite wild... Overfitting?

k -Nearest neighbors

Get more "comps"

More reliable estimate if you base estimate off of a larger set of comparable homes

↑ \$ = ???

↑ \$ = 850k

↑ \$ = 833k

↑ \$ = 749k

↑ \$ = 901k

k-NN regression more formally

Dataset of (Int , $\$$) pairs: $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$

Query point: x_q

1. Find k closest x_i in dataset

$(x_{NN_1}, x_{NN_2}, \dots, x_{NN_k})$ such that for any x_i not in nearest neighbor set, $\text{distance}(x_i, x_q) \geq \text{distance}(x_{NN_k}, x_q)$

2. Predict

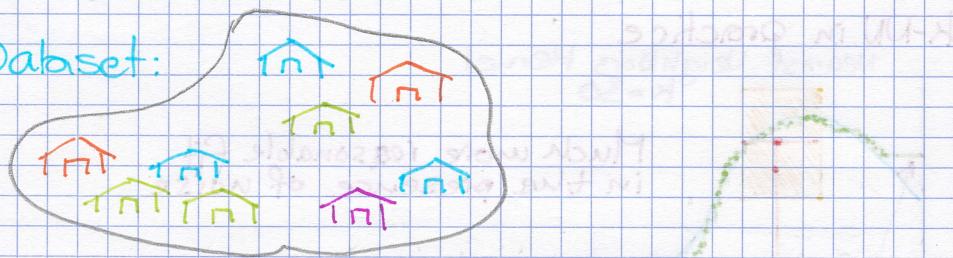
$$\hat{y}_q = \frac{1}{k} (y_{NN_1} + y_{NN_2} + \dots + y_{NN_k})$$

$$= \frac{1}{k} \sum_{j=1}^k y_{NN_j}$$

Performing k-NN search

- Query house: Int

- Dataset:



- Specify: Distance metric

- Output: Most similar houses



k -NN algorithm

Initialize $\text{Dist2kNN} = \text{sort}(\delta_1, \dots, \delta_k) \leftarrow$ list of sorted distances

Sort first k houses

by distance to query house

For $i = k+1, \dots, N$

Compute $\delta = \text{distance}(\hat{\text{h}}_i, \hat{\text{h}}_{\text{q}})$ query house

if $\delta < \text{Dist2kNN}[k]$

find j such that $\delta > \text{Dist2kNN}[j-1]$
but $\delta < \text{Dist2kNN}[j]$ remove
furthest house and shift queue:

$$\hat{\text{h}}[j+1:k] = \hat{\text{h}}[j:k-1]$$

$$\text{Dist2kNN}[j+1:k] = \text{Dist2kNN}[j:k-1]$$

set $\text{Dist2kNN}[j] = \delta$ and $\hat{\text{h}}[j] = \hat{\text{h}}_i$

Return k most similar houses $\hat{\text{h}}$

closest houses to query house $\hat{\text{h}}_1 \hat{\text{h}}_2 \hat{\text{h}}_3 \hat{\text{h}}_4$

k -NN in practice

Nearest Neighbors Kernel
 $k=30$



Much more reasonable fit
in the presence of noise

Discontinuities!
Neighbor either in or out

Boundary & sparse region issues

Issues with discontinuities

Overall predictive accuracy might be okay, but...

For example, in housing application:

- If you are a buyer or seller, this matters
- Can be a jump in estimated value of house going from 2640 sq.ft to 2641 sq.ft.
- Don't really believe this type of fit.

Weighted k-nearest neighbors

Weigh more similar houses more than those less similar in list of k-NN

Predict: weights on NN

$$\hat{y}_q = \frac{c_{qNN1}y_{NN1} + c_{qNN2}y_{NN2} + c_{qNN3}y_{NN3} + \dots + c_{qNNK}y_{NNK}}{\sum_{j=1}^k c_{qNNj}}$$

How to define weights?

Want weight c_{qNNj} to be small when distance (x_{NNj}, x_q) large

and c_{qNNj} to be large when distance (x_{NNj}, x_q) small

Simple method:

$$c_{qNNj} = \frac{1}{\text{distance}(x_j, x_q)}$$

Kernel weights for $d=1$

Define: $c_{qNNj} = \text{Kernel}_\lambda(|x_{NNj} - x_q|)$ Simple isotropic case

Gaussian kernel:

$$\text{Kernel}_\lambda(|x_i - x_q|) = \exp(-(|x_i - x_q|)^2 / \lambda)$$

Note: never exactly 0!

Kernel weights for $d \geq 1$

Define: $c_{qNNj} = \text{Kernel}_\lambda(\text{distance}(x_{NNj}, x_q))$

Kernel regression

Instead of just weighting NN, weight all points

Predict weight on each datapoint

$$\hat{y}_q = \frac{\sum_{i=1}^N c_{qi} y_i}{\sum_{i=1}^N c_{qi}} = \frac{\sum_{i=1}^N \text{Kernel}_\lambda(\text{distance}(x_i, x_q)) * y_i}{\sum_{i=1}^N \text{Kernel}_\lambda(\text{distance}(x_i, x_q))}$$

| Nadaraya-Watson
| Kernel weighted average |

Epanechnikov Kernel ($\lambda=0, 2$)



Kernel has bounded support... Only subset of data needed to compute local fit

Choice of bandwidth λ

Often choice of kernel matters much less than choice of λ

Choosing λ (or k in k-NN)

How to choose? Same story as always...

Cross Validation

Formalizing the idea of local fits

Contrasting with global average

A globally constant fit weights all points equally

$$y_p = \frac{1}{N} \sum_{i=1}^N y_i = \frac{\sum_{i=1}^N c y_i}{\sum_{i=1}^N c}$$

equal weight on each datapoint

Kernel regression leads to locally constant fit
- slowly add in some points and let others gradually die off

$$y_p = \frac{\sum_{i=1}^N \text{Kernel}_X(\text{distance}(x_i, x_q)) * y_i}{\sum_{i=1}^N \text{Kernel}_X(\text{distance}(x_i, x_q))}$$

Local linear regression

So far, discussed fitting constant function locally at each point

→ "locally weighted averages"

Can instead fit a line or polynomial locally at each point

→ "locally weighted linear regression"

Local regression rules of thumb

- Local linear fit reduces bias at boundaries with minimum increase in variance
- Local quadratic fit doesn't help at boundaries and increases variance, but does help capture curvature in the interior
- With sufficient data, local polynomials of odd degree dominate those of even degree

Recommended default choice:
Local linear regression

Discussion on k-NN and kernel regression

Nonparametric approaches

k-NN and kernel regression are examples of nonparametric regression

General goals of nonparametrics:

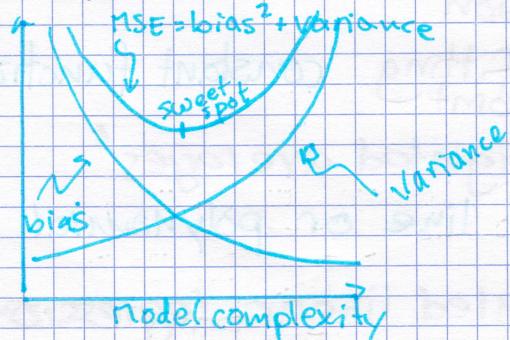
- Flexibility
- Make few assumptions about $f(x)$
- Complexity can grow with the number of observations N

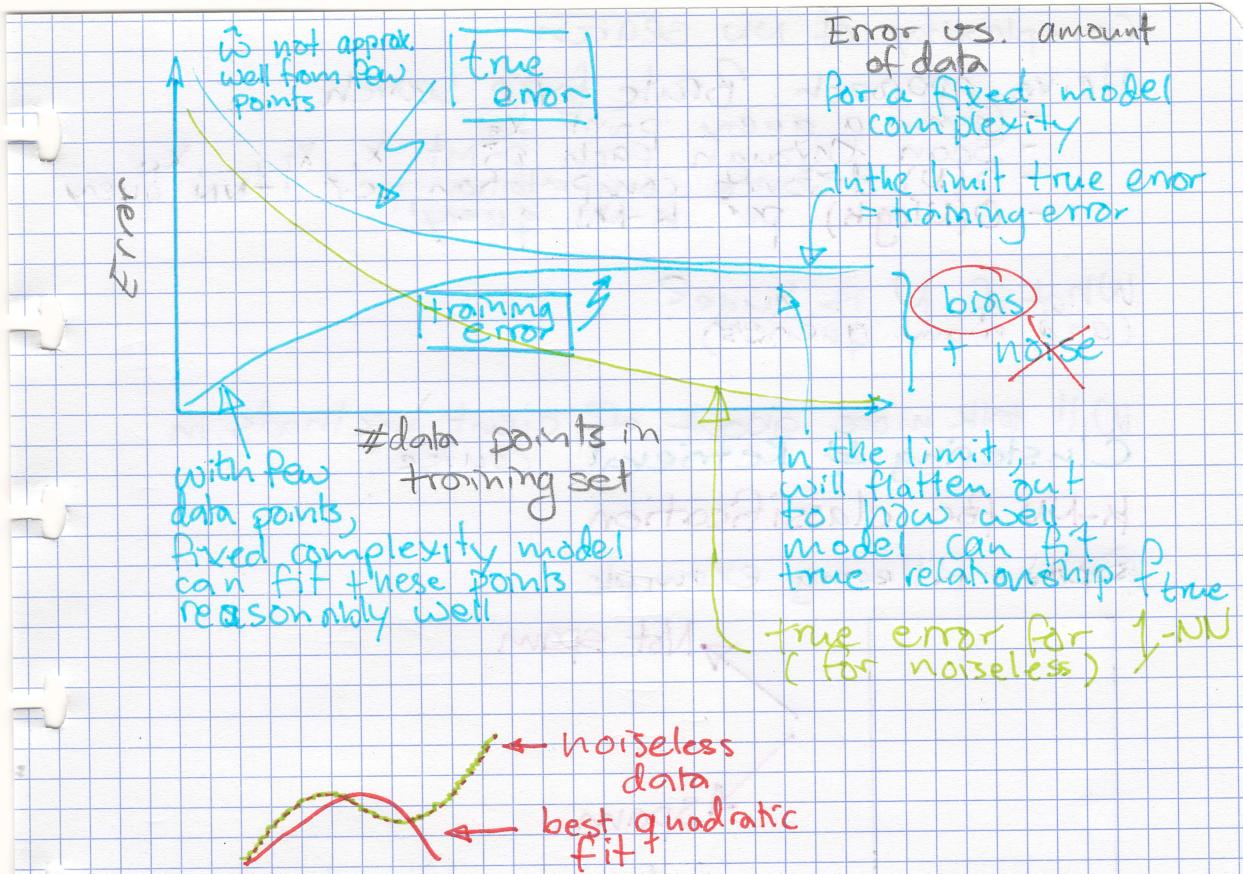
Lots of other choices:

- Splines, trees, locally weighted structured regression models.

Limiting behavior of NN:
Noiseless setting $\epsilon_i = 0$

In the limit of getting an infinite amount of noiseless data, the MSE of 1-NN fit goes to 0





Limiting behavior of NN:
Noisy data setting

In the limit of getting an infinite amount of data, the MSE of NN fit goes to 0 if k grows, too

Issues with high-dimensions, data scarcity, and computational complexity

NN and Kernel methods for large d or small N

NN and Kernel methods work well when the data cover the space, but:

- the more dimensions d you have, the more points N you need to cover the space
- need $N = O(\exp(d))$ data points for good performance

This is where parametric models become useful...

Complexity of NN search

Naive approach. Brute force search

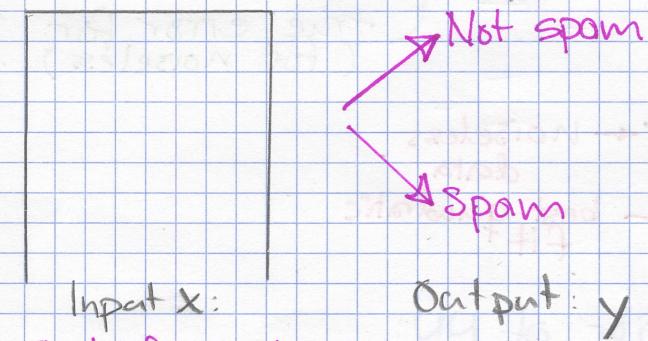
- Given a query point x_q
- Scan through each point x_1, x_2, \dots, x_N
- $O(N)$ distance computation per 1-NN query!
- $O(N \log k)$ per k-NN query!

What if N is huge?
(and many queries)

Will talk more about efficient methods in
Clustering & Retrieval course

k-NN for classification

Spam filtering example



Text of email
sender, IP, ...

Using k-NN for classification

Space of labeled emails (not spam vs. spam),
organized by similarity of text

