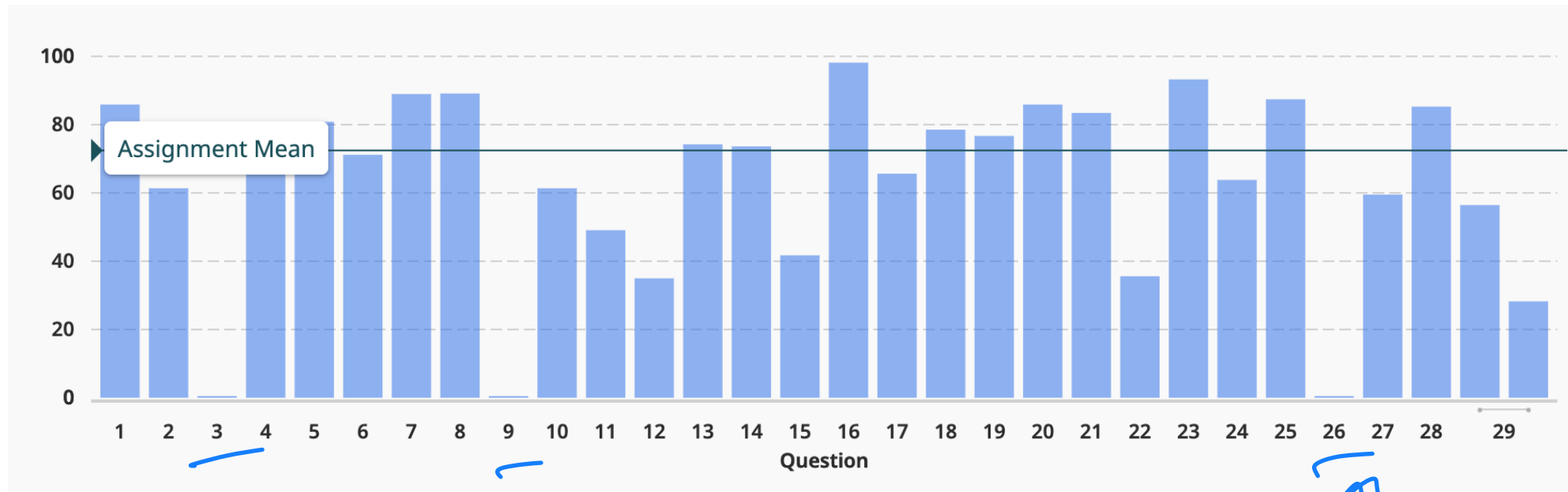


Midterm grades adjusted

- 3 hardest questions became bonus questions



Midterm Exam 27.0 points

Minimum

34.26%

Median

73.15%

Maximum

108.33%

Mean

72.22%

Std Dev ?

16.14%

Non-parametric methods

Nearest Neighbours

Natasha Jaques

Parametric vs non-parametric

k NN

- A model is parametric if # parameters does not depend on # samples

OLS = linear regression
neural networks

// learning a bunch of parameters / weights

- A model is non-parametric if # parameters increases with # samples

- Does not mean absence of parameters!

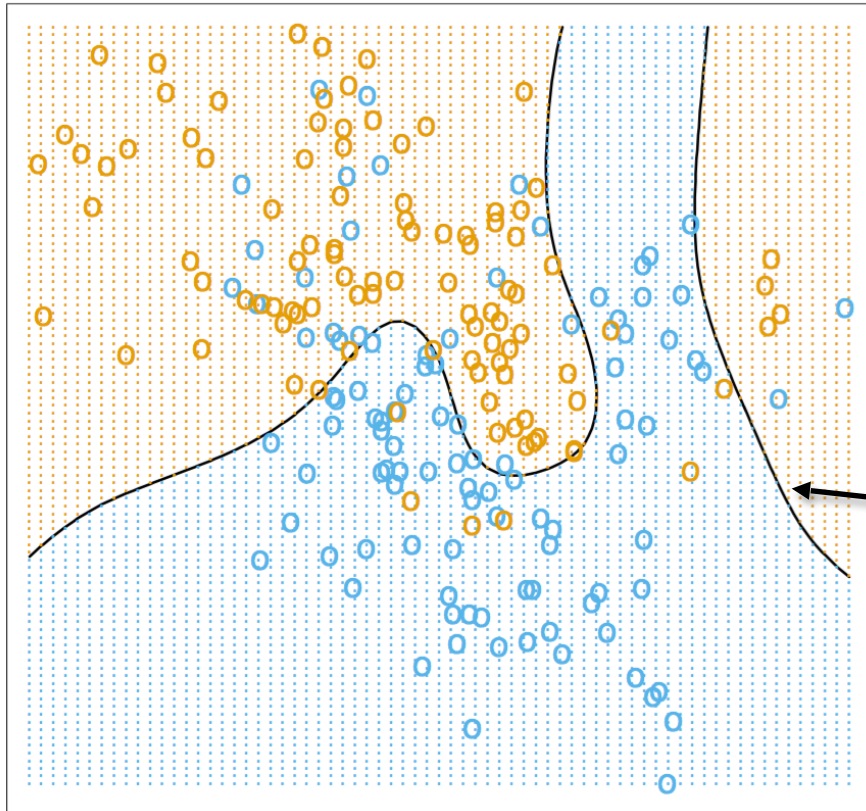
→ today's class

This lecture: k nearest neighbors

simple

- Assume we have a classification task
- To classify a new point x :
 - Find its k nearest neighbors in the training data → smallest distance in feature space (2 features)
 - Set y to be the majority vote of the labels of these nearest neighbors
- Design choices / hyperparameters:
 - Number of nearest neighbors k
 - Distance metric
 - Aggregation method

Example: Bayes classifier



Training data:

- True label: +1
- True label: -1

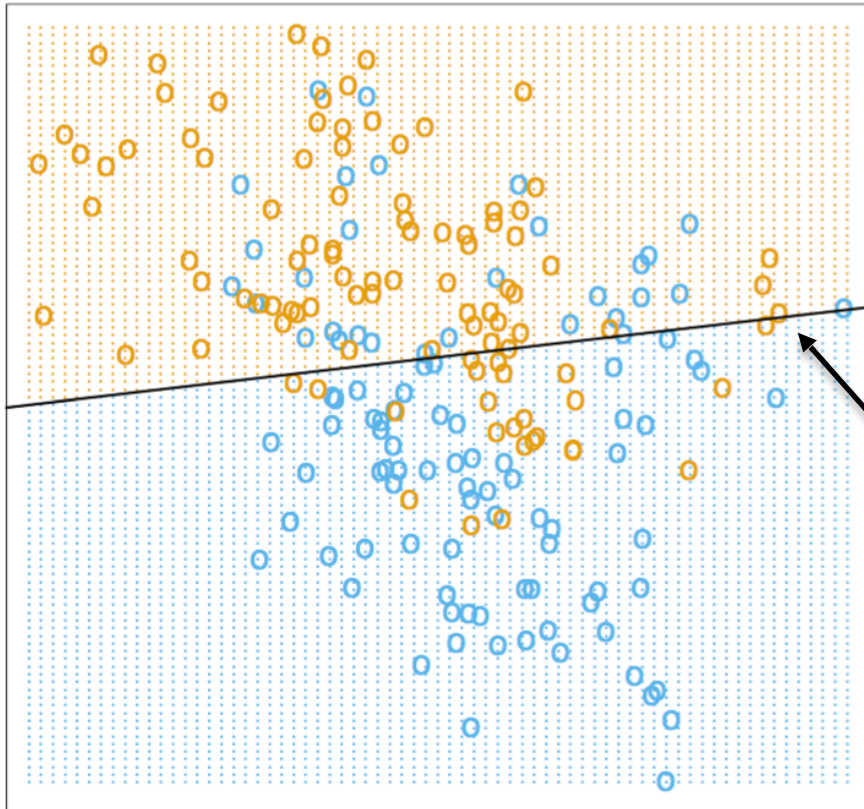
Optimal Bayes classifier:

$$\mathbb{P}(Y = 1 | X = x) = \frac{1}{2}$$

← misclassification

- ▢ Predicted label: +1
- ▢ Predicted label: -1

Linear decision boundary



Training data:

○ True label: +1

○ True label: -1

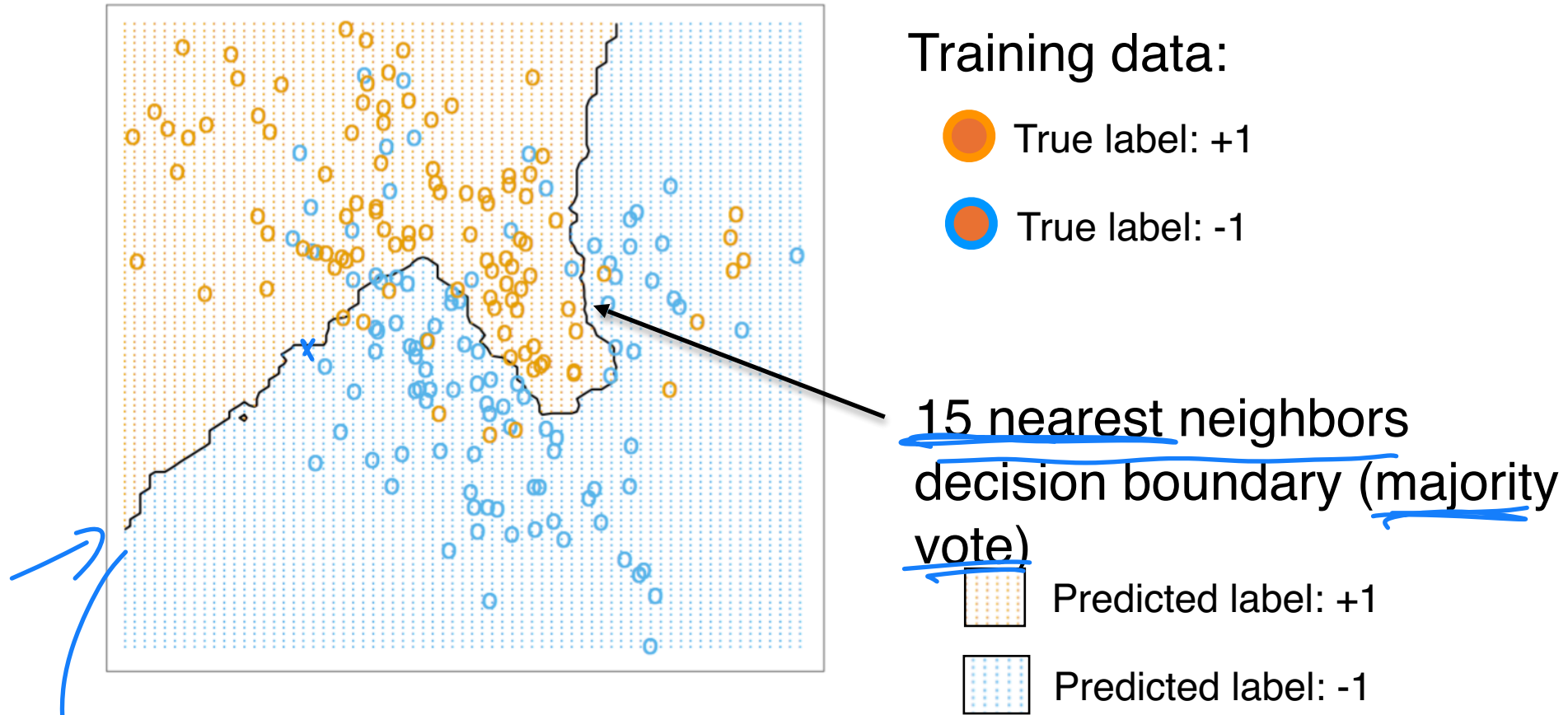
Learned linear decision boundary:

$$x^T w + b = 0$$

▨ Predicted label: +1

▨ Predicted label: -1

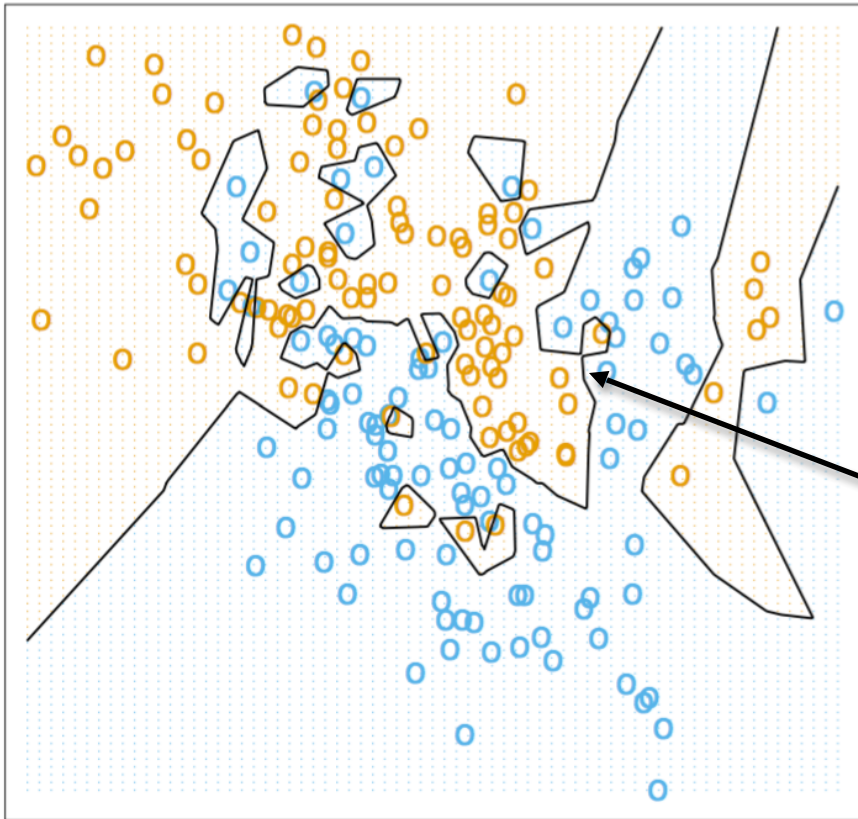
$k = 15$ nearest neighbors boundary



boundary depends on
the training data points

$k = 1$ nearest neighbor boundary

→ overfitting



Training data:

○ True label: +1

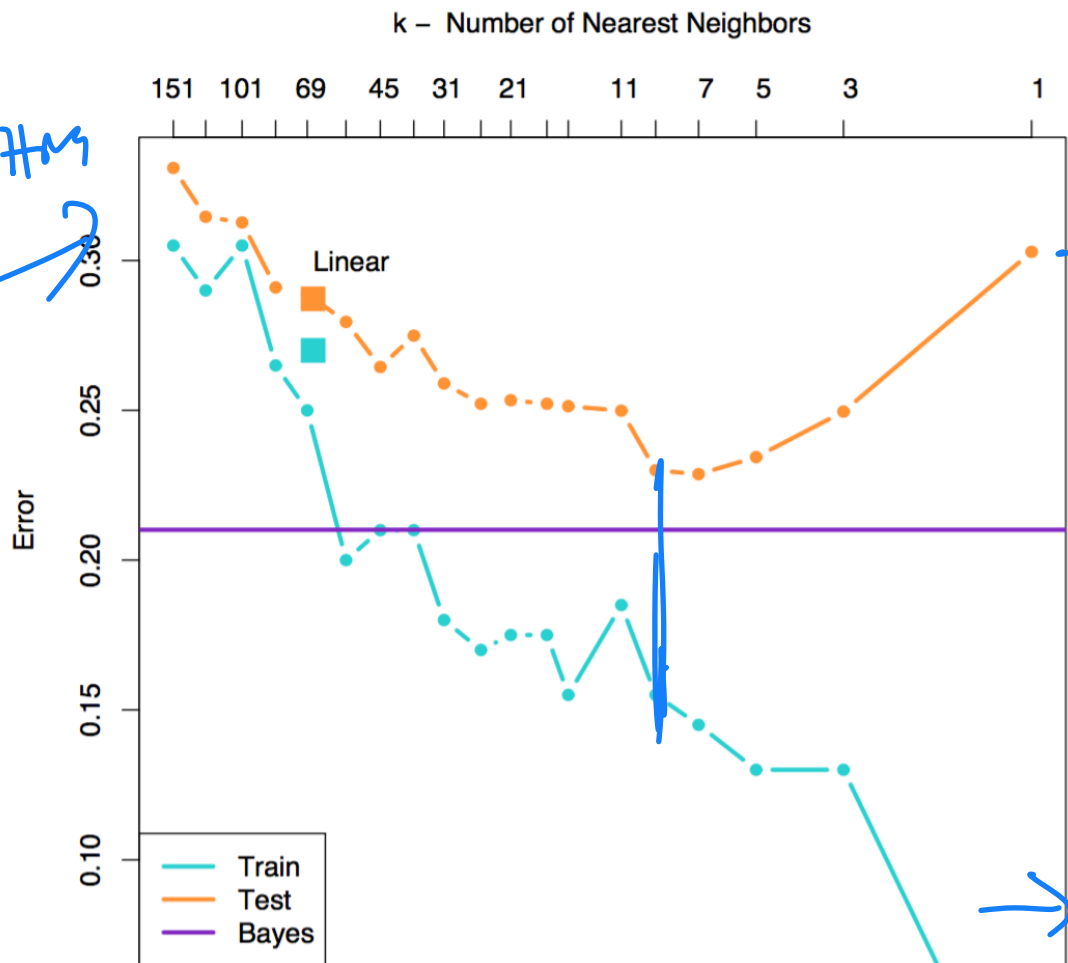
○ True label: -1

1 nearest neighbor decision boundary (majority vote)

■ Predicted label: +1

■ Predicted label: -1

k nearest neighbors error



Parametric vs non-parametric

After training, discard your data

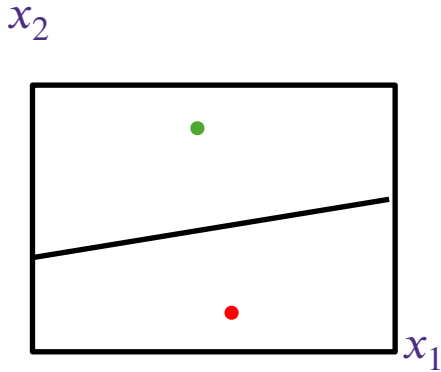
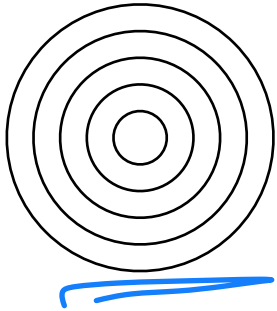
↳ keep training data around for inference
eg. KNN

Interpretable

Notable distance metrics & level sets

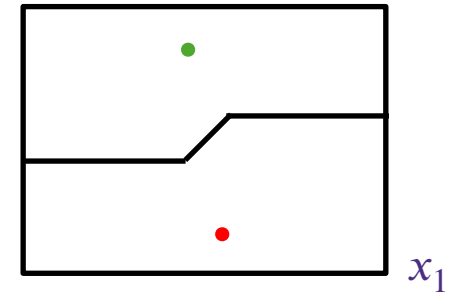
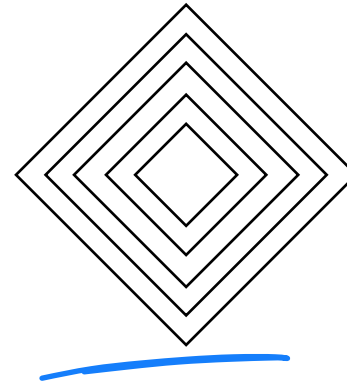
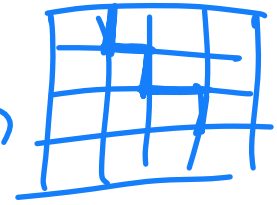
ℓ_2 norm (Euclidean)

$$d(u, v) = \|u - v\|^2$$



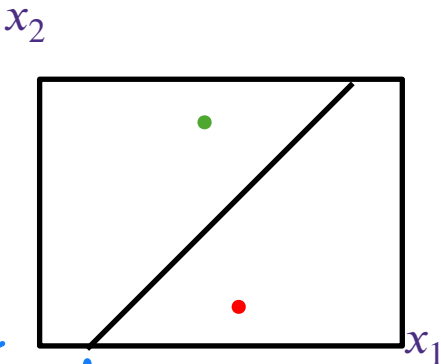
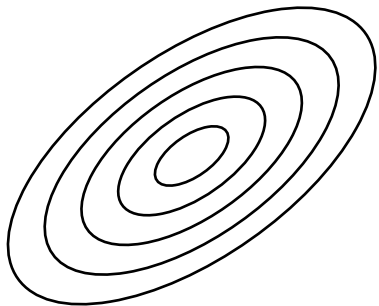
ℓ_1 norm (Manhattan, taxicab)

$$\|u - v\|_1 = \sum_i |u_i - v_i|$$



Mahalanobis norm

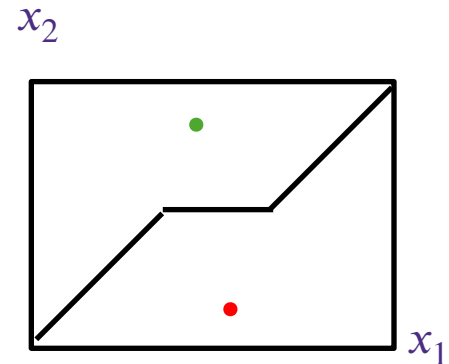
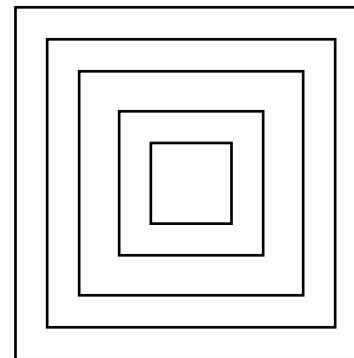
$$(u - v)^T M (u - v)$$



weight diff dimensions differently

ℓ_∞ norm (max)

\rightarrow max dist b/w vectors

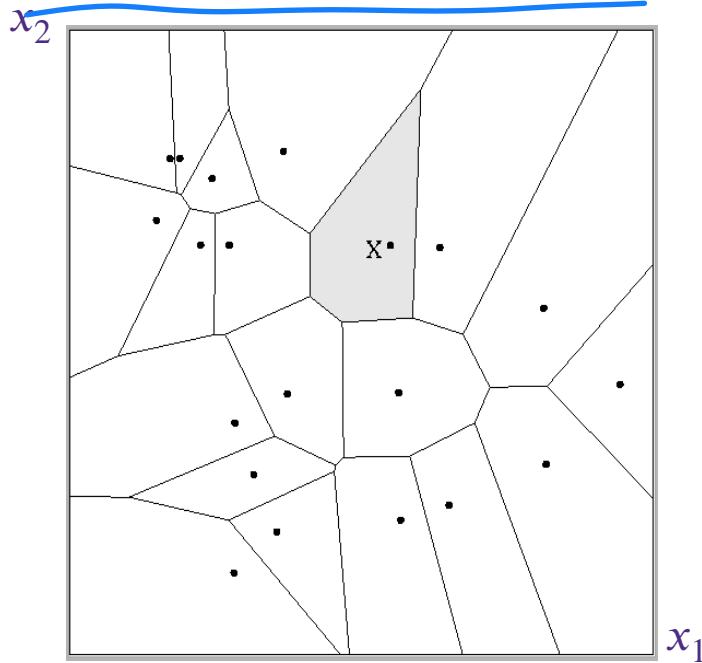


Example: distance metrics with $k = 1$ NN

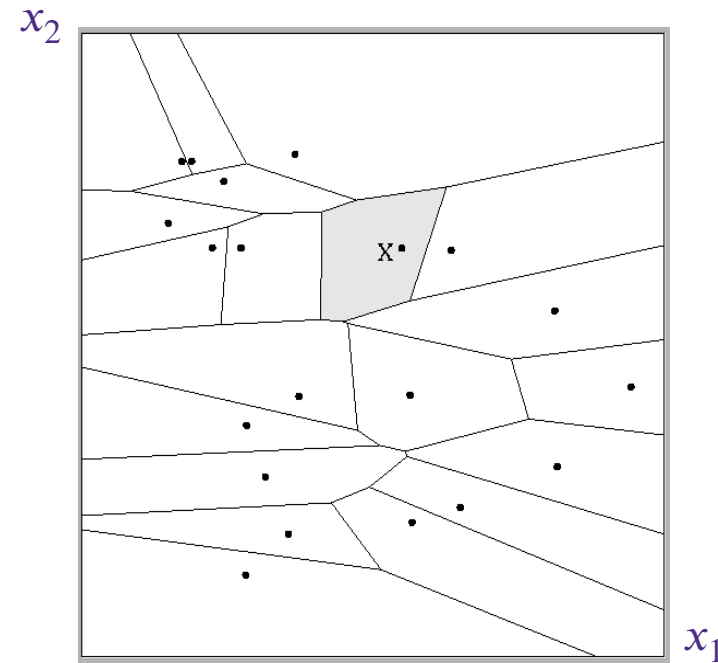
12

Mahalanobis

$$d(x, x') = (x_1 - x'_1)^2 + (x_2 - x'_2)^2$$



$$d(x, x') = (x_1 - x'_1)^2 + 9(x_2 - x'_2)^2$$



// know geometry of feat space

Learned distance metrics

Training data



Dog



Cat

Test data



learn distance with nn

Cool online kNN demo (credit to Vicky Ye)

<http://vision.stanford.edu/teaching/cs231n-demos/knn/>

→ 1-NN classification: Theoretical guarantees

→ have enough data
 → function is smooth

$$D = \{ (x^{(i)}, y^{(i)}) \}_{i=1}^n \sim P \quad x^{(i)} \in \mathbb{R}^d, y \in \{0, 1\}$$

Given test point x , let x_{NN} be the nearest neighbour in D

Error if $y_{NN} \neq y$

Case 1: $y_{NN} = 1, y = 0$ w.p. $P(y=1 | x_{NN}) P(y=0 | x)$

Case 2: $y_{NN} = 0, y = 1$ w.p. $P(y=0 | x_{NN}) P(y=1 | x)$

As $n \rightarrow \infty$, $P(y | x_{NN}) \rightarrow P(y | x)$ // distance b/w x & $x_{NN} \rightarrow \underline{0}$
 as $n \rightarrow \infty$

Error: $2 P(y=1 | x) P(y=0 | x)$ as $n \rightarrow \infty$

$$= 2 p^* (1 - p^*)$$

$$\leq 2 p^*$$

Bayes error

define $p^* = \min \{ P(y=1 | x), P(y=0 | x) \}$

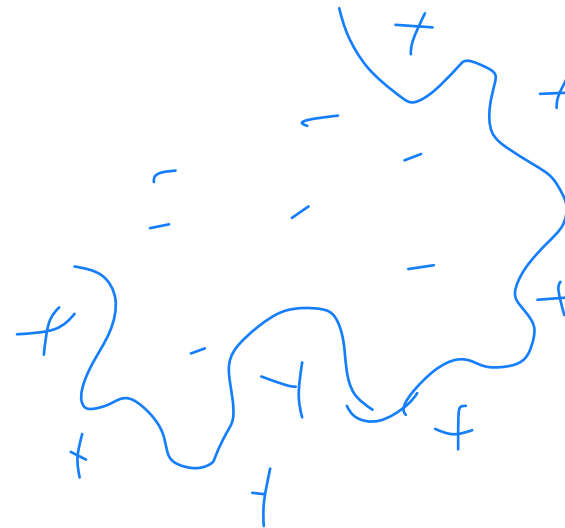
best error I can get measurement error

↑ with smoothness

1-NN classification: Theoretical guarantees

Theorem[Cover, Hart, 1967] If P_X is supported everywhere in \mathbb{R}^d and $P(Y = 1|X = x)$ is smooth everywhere, then as $n \rightarrow \infty$ the 1-NN classification rule has error at most twice the Bayes error rate.

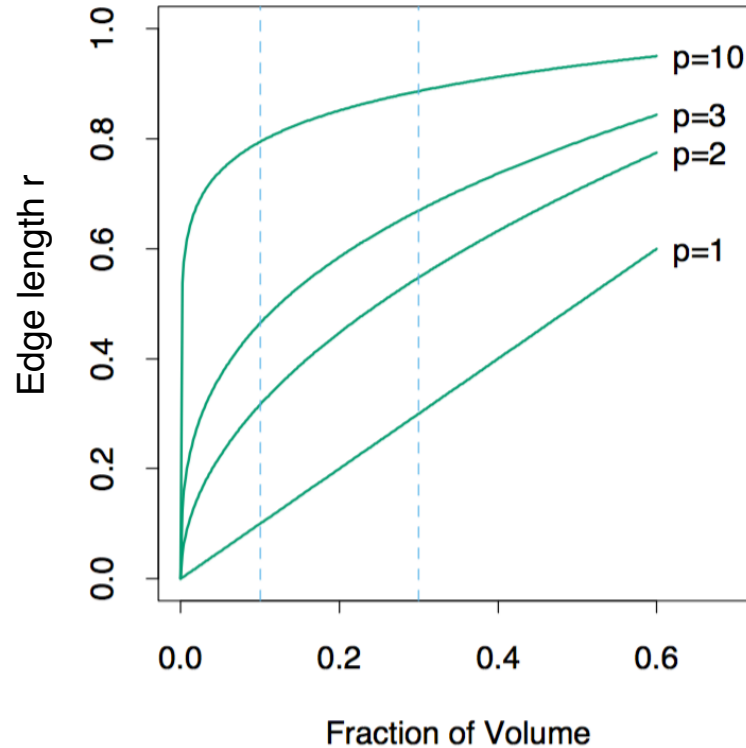
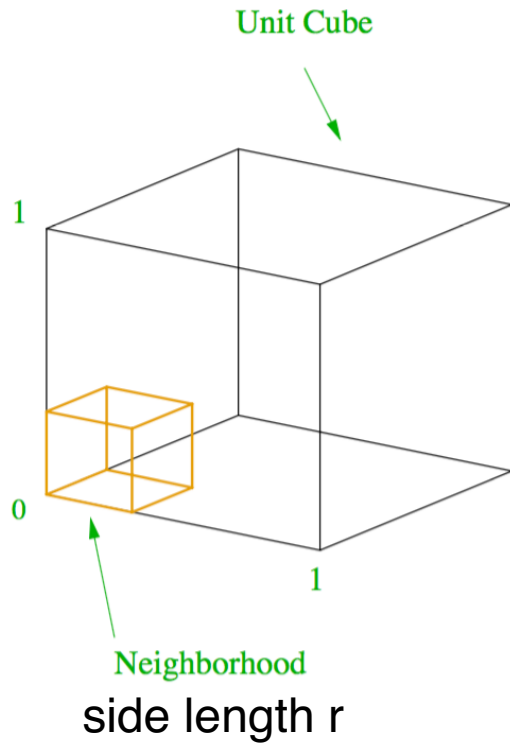
Can fit arbitrarily
complicated
decision boundaries



Curse of dimensionality, example 1

of samples in D we would need to have to find a neighbour with size distance

↓ increases exponentially!



$$0 \text{---} r \text{---} 1$$

$$0.3 = r = p$$

$$0.3^2 = 0.09$$

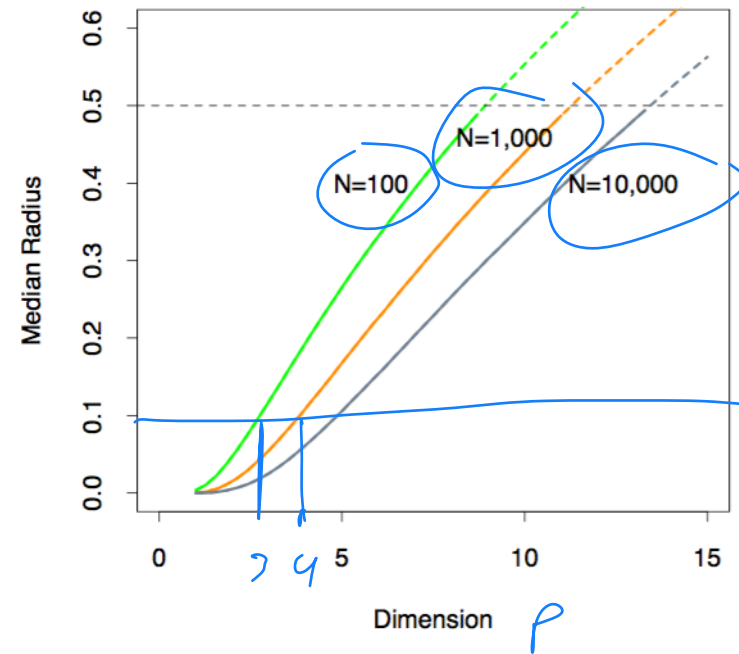
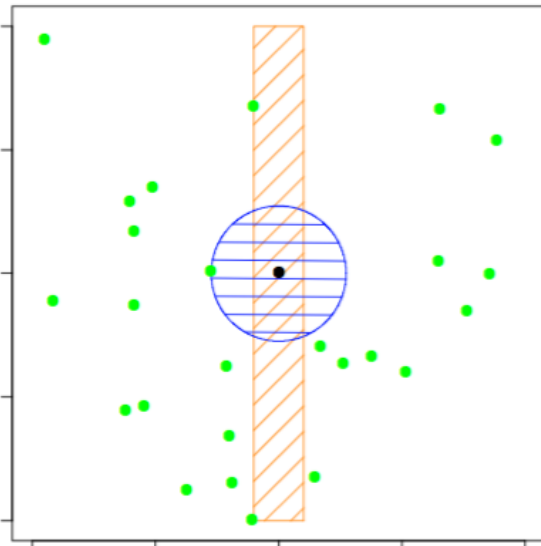
X is uniformly distributed over $[0, 1]^p$. What is $\mathbb{P}(X \in [0, r]^p)$? $= \frac{1}{r^p}$

How many samples do we need so that a nearest neighbor is within a cube of side length r ?

Curse of dimensionality, example 2

high d
space
is hard
to cover

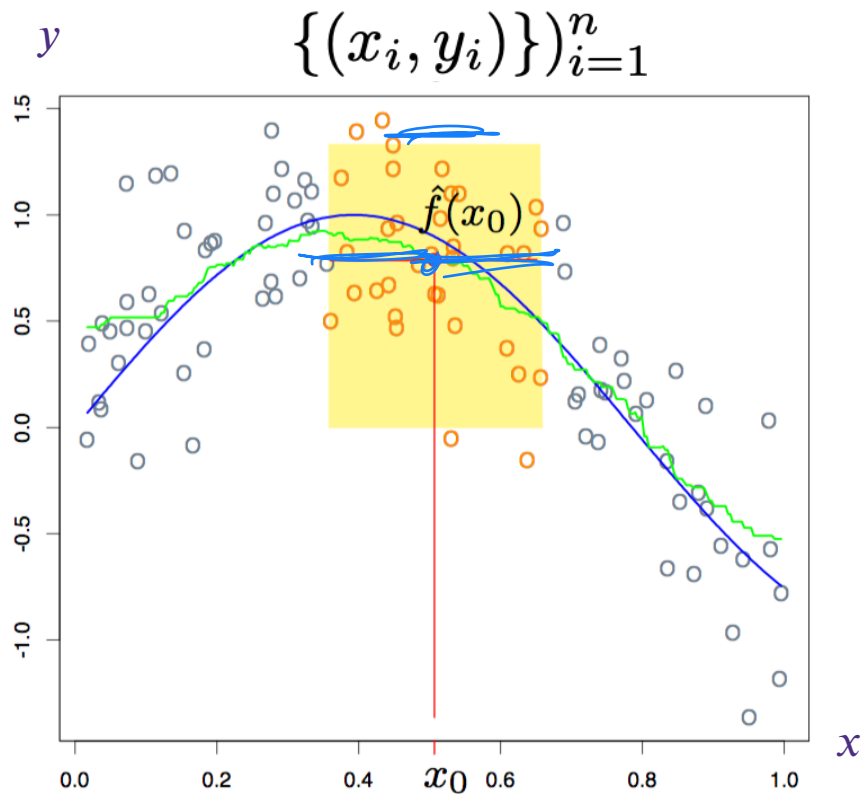
$\{X_i\}_{i=1}^n$ are uniformly distributed over $[-.5, .5]^p$.



What is the median distance from a point at origin to its 1NN?

How many samples do we need so that a median Euclidean distance is within r ?

Nearest neighbor regression



k-NN regressor

$$\hat{f}(x) = \frac{1}{k} \sum_{i \in k\text{NN}(x)} y^{(i)}$$

$$= \sum_{i=1}^n y^{(i)} \mathbb{1}_{\{x^{(i)} \in k\text{NN}(x)\}}$$

recall

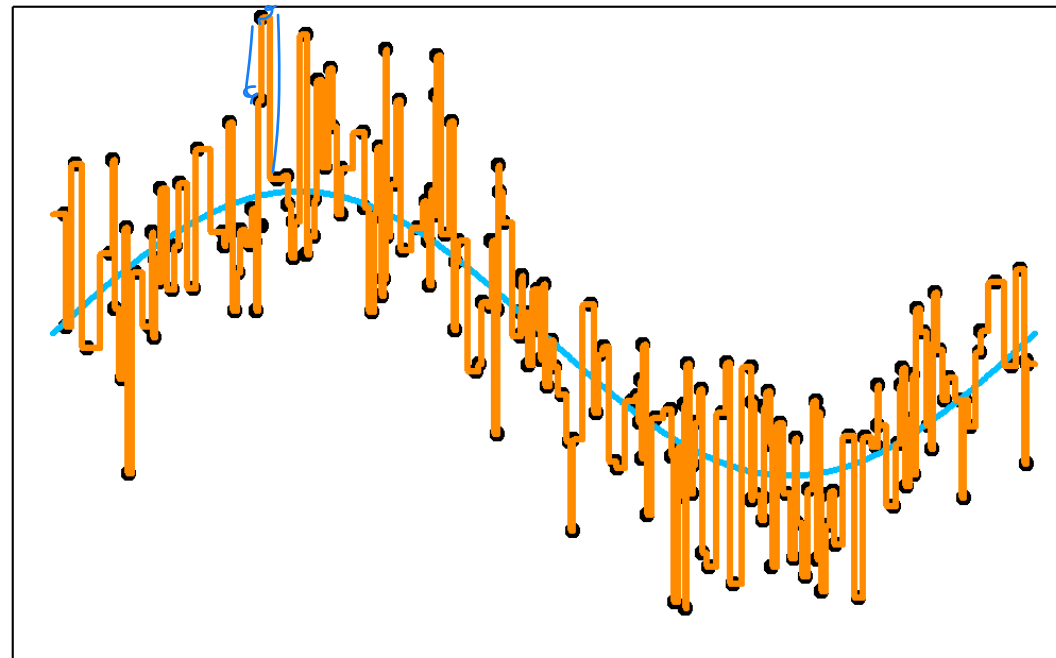
$$f^*(x) = \arg \min_f \mathbb{E}[(f(x) - y)^2]$$

$$= \mathbb{E}[y | x]$$

// best in theory

Overfitting

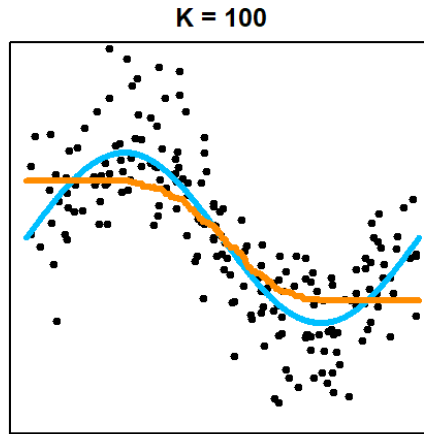
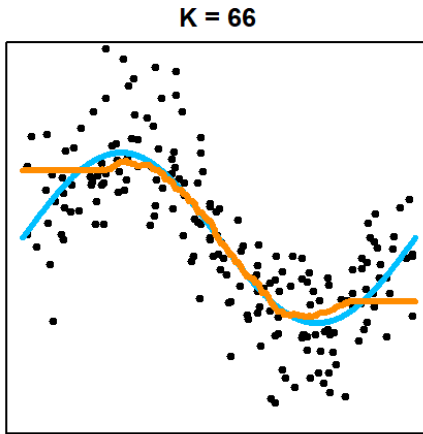
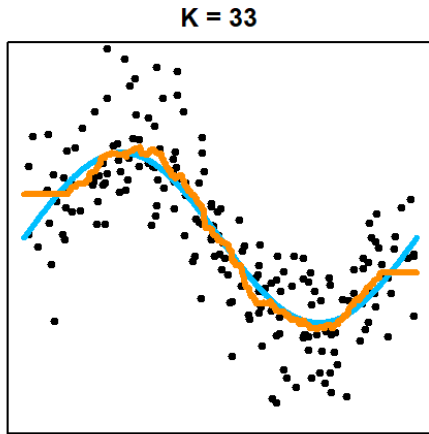
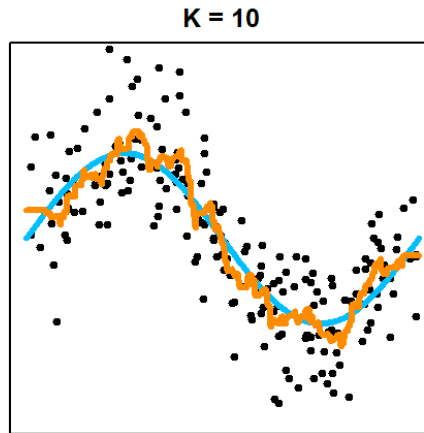
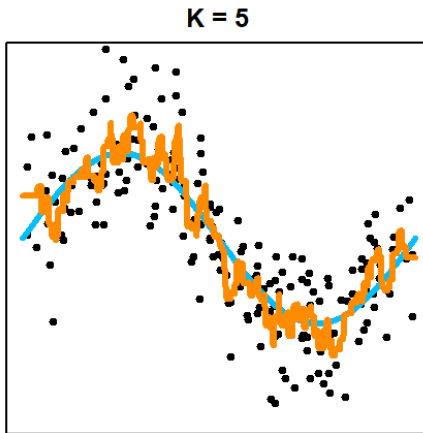
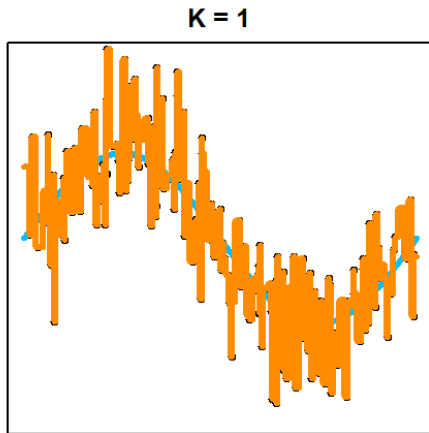
1-Nearest Neighbor Regression



Bias vs variance

(A)

highest variance



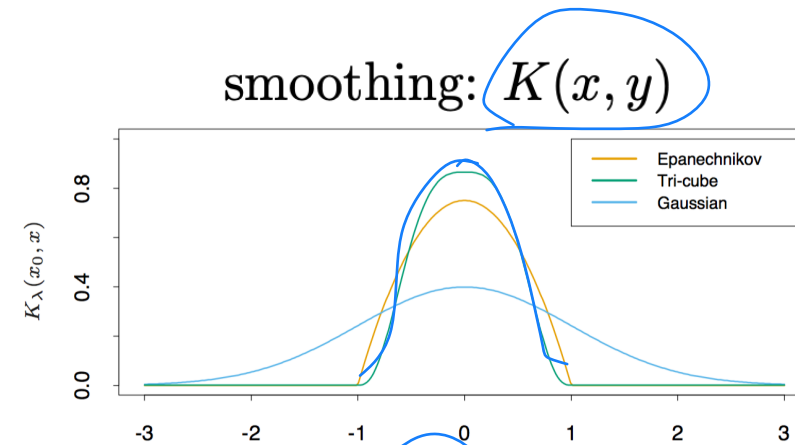
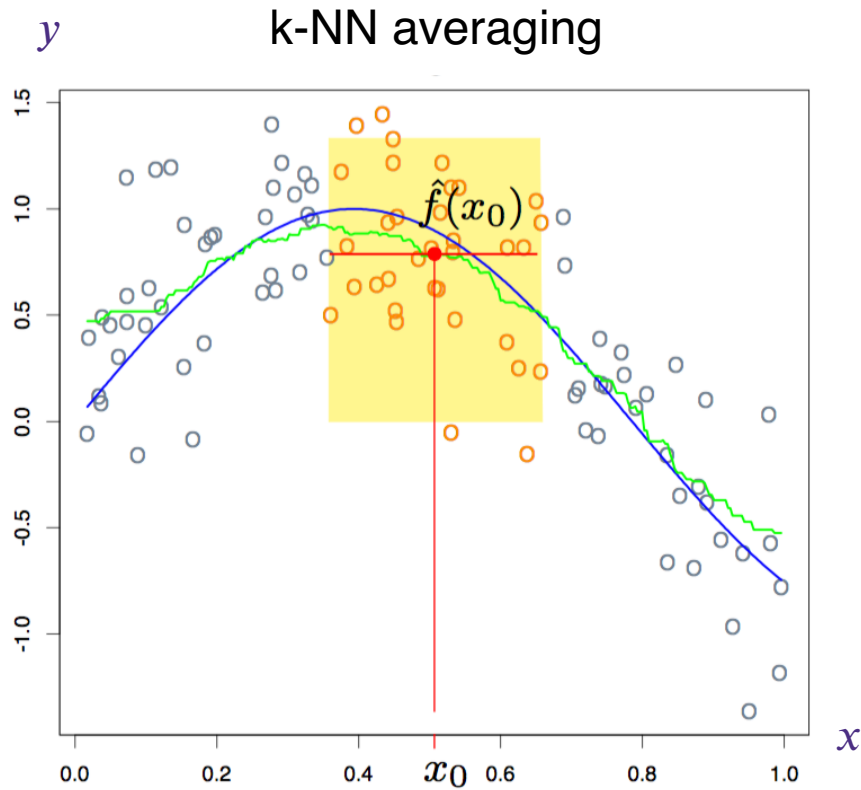
high bias

D

E

F

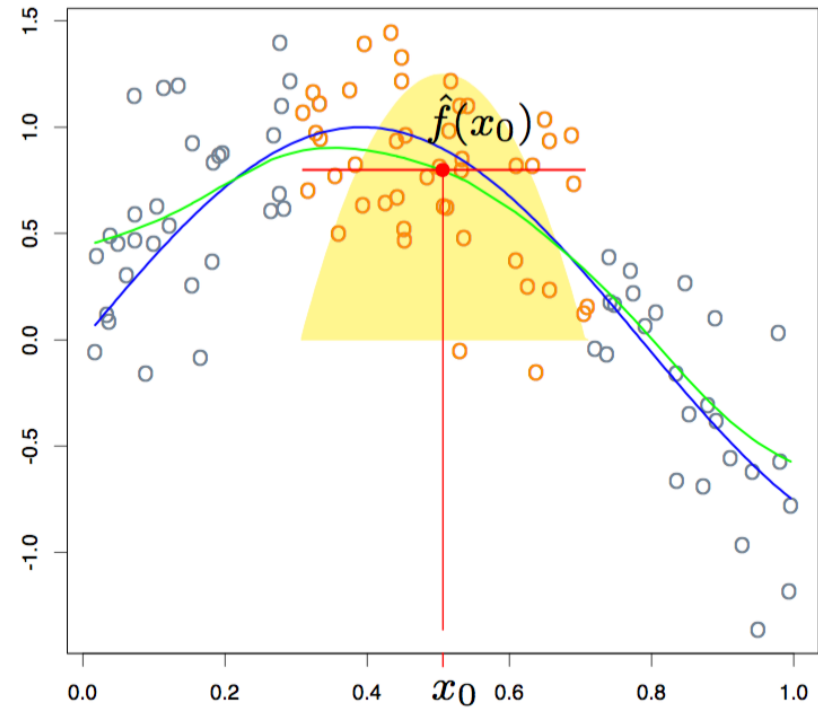
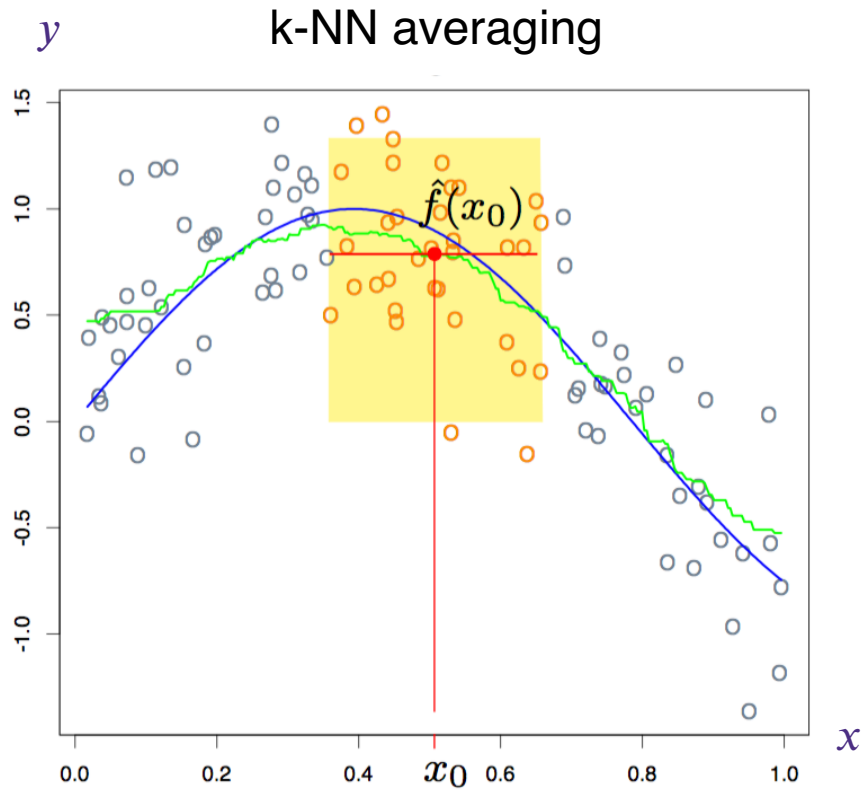
Smoothed nearest neighbor regression



$$\hat{f}(x) = \frac{\sum_{i=1}^n K(x, x^{(i)}) y^{(i)}}{\sum_{i=1}^n K(x, x^{(i)})}$$

Smoothed nearest neighbor regression

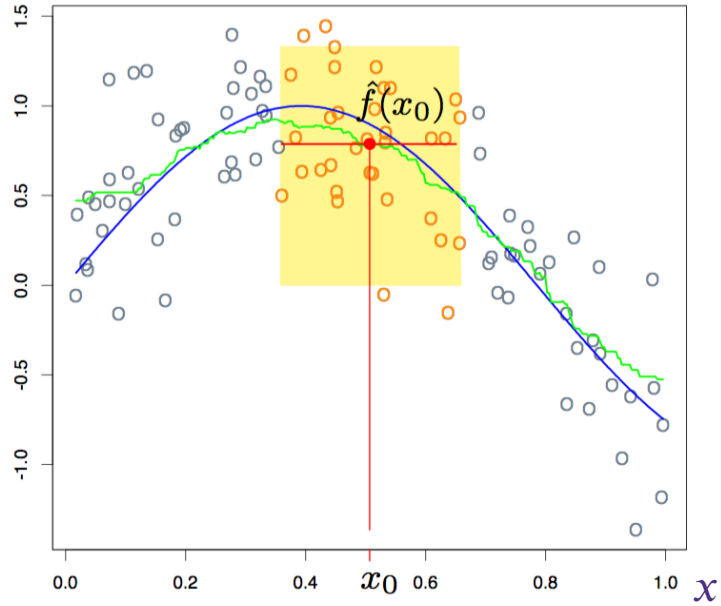
$$\hat{f}(x) = \frac{\sum_{i=1}^n K(x, x^{(i)}) y^{(i)}}{\sum_{i=1}^n K(x, x^{(i)})}$$



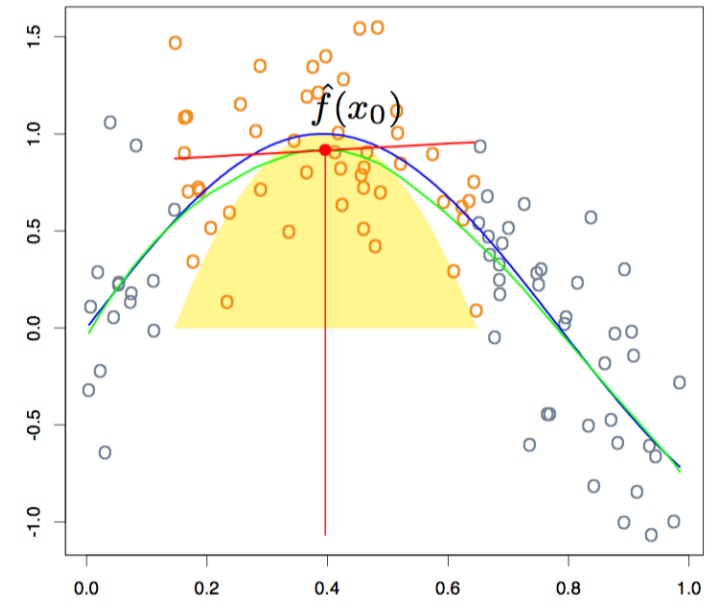
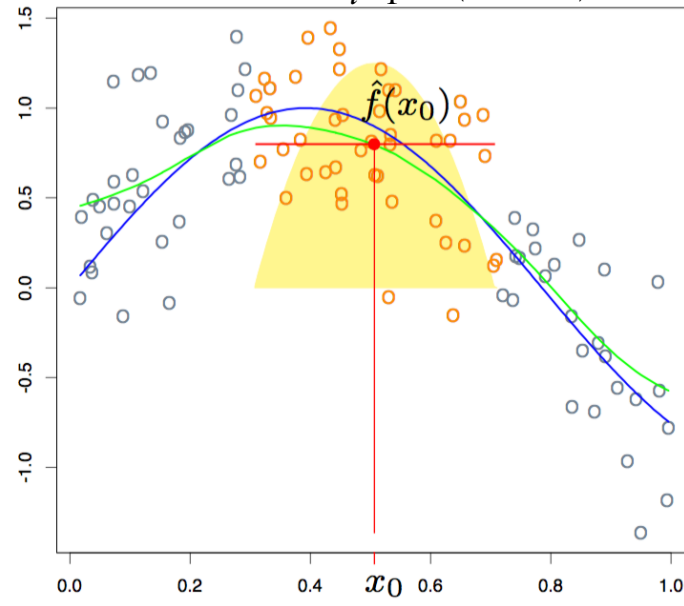
Locally linear regression

y

k-NN averaging



$$\hat{f}(x) = \frac{\sum_{i=1}^n K(x, x^{(i)}) y^{(i)}}{\sum_{i=1}^n K(x, x^{(i)})}$$



Have we seen non-parametric methods before?

- Kernel methods are non-parametric:

$$f(x) = \sum_{i=1}^n \alpha_i K(x^{(i)}, x)$$

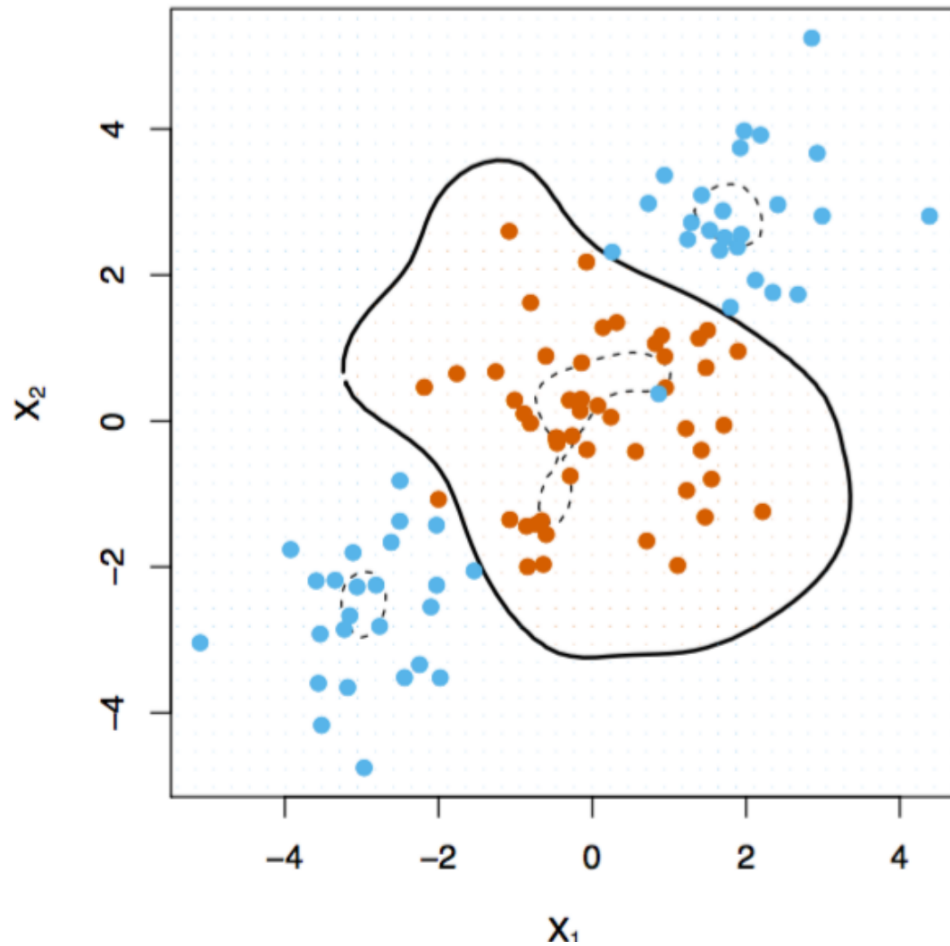
parameters goes up with # data

- Compare with (smoothed) nearest neighbors:

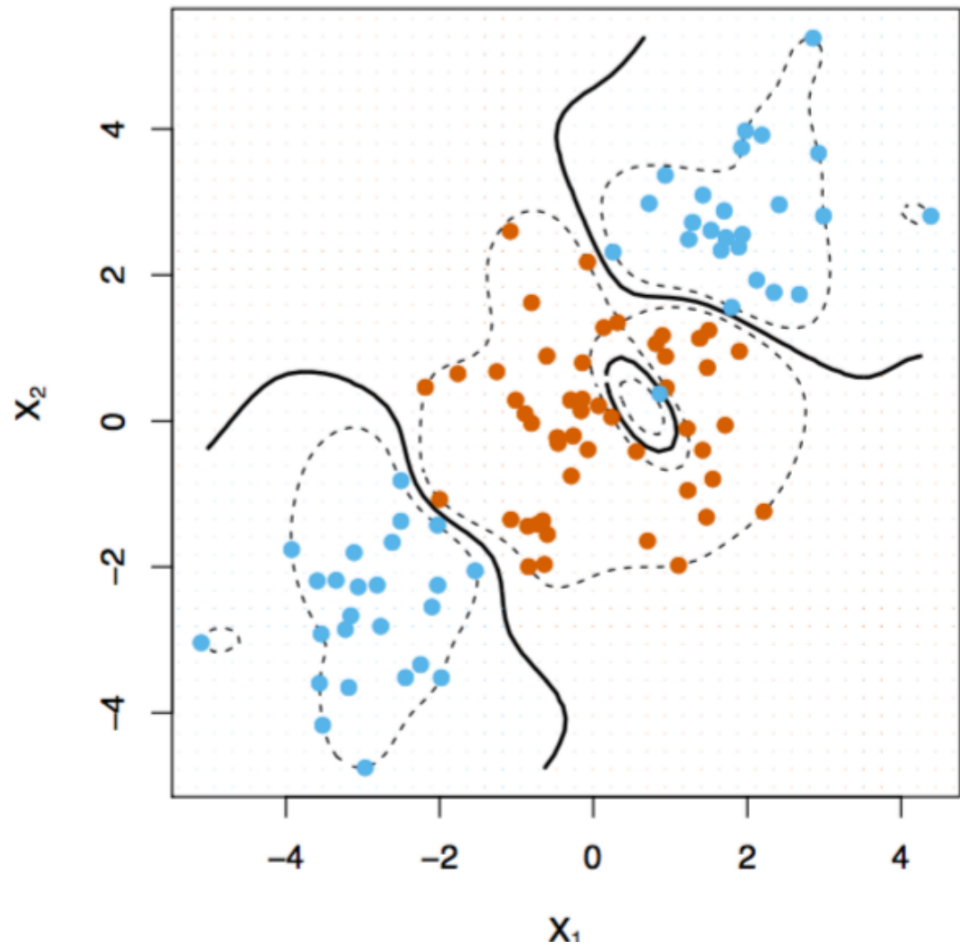
$$\hat{f}(x) = \frac{\sum_{i=1}^n K(x, x^{(i)}) y^{(i)}}{\sum_{i=1}^n K(x, x^{(i)})}$$

The Radial Basis Function (RBF) kernel $\exp\left(-\frac{\|x^{(i)} - x\|_2^2}{2\sigma^2}\right)$

Bandwidth σ is large enough

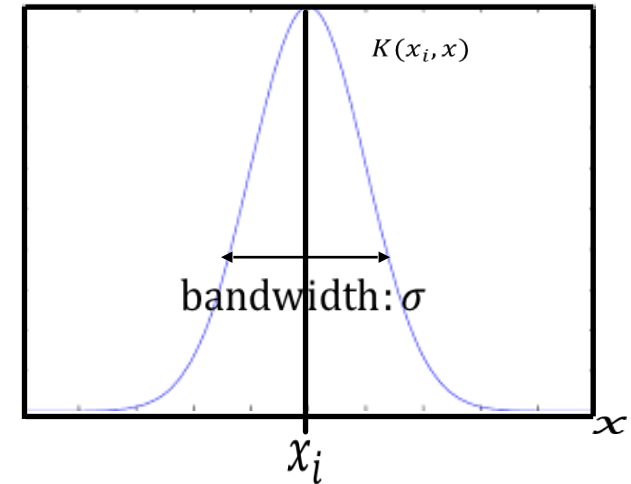
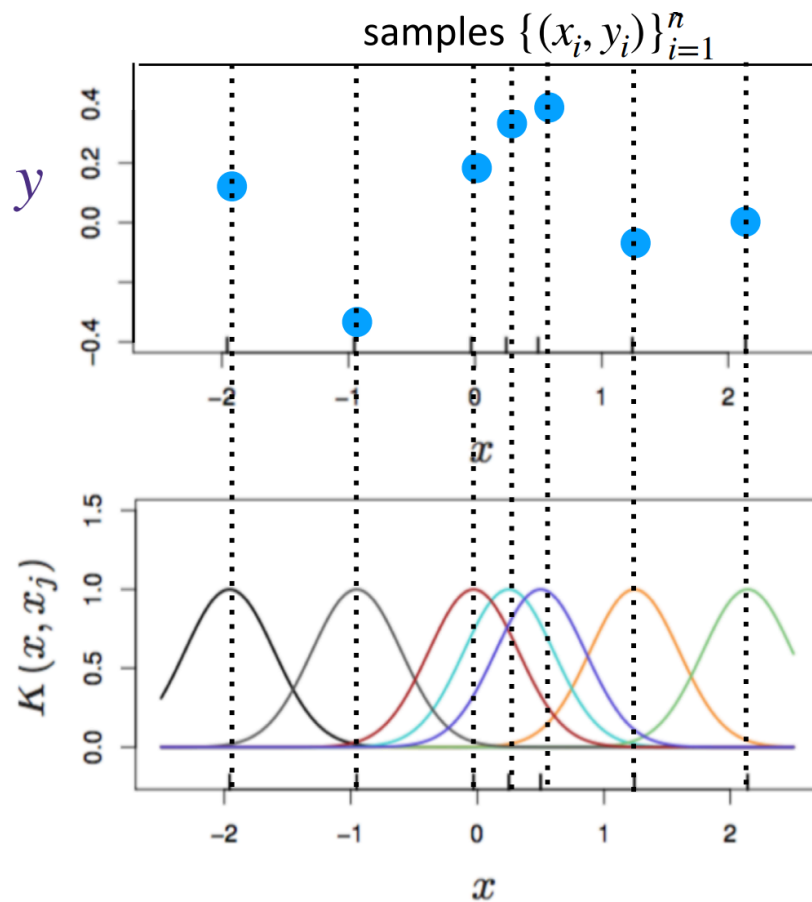


Bandwidth σ is small



The Radial Basis Function (RBF) kernel $\exp\left(-\frac{\|x^{(i)} - x\|_2^2}{2\sigma^2}\right)$

$$f(x) = \sum_{i=1}^n \alpha_i K(x^{(i)}, x)$$



Kernel methods

$$f(x) = \sum_{i=1}^n \alpha_i K(x^{(i)}, x)$$

- Can be seen as a soft, learned version of “nearest” neighbors
- $K(x^{(i)}, x) = \phi(x^{(i)})^\top \phi(x)$ defines “similarity” between $x^{(i)}$ and x
- How many parameters?

→ n

Takeaways

- k-NN is very simple to explain and implement
- No training! But inference can still be computationally demanding.
- You can use other forms of distance (not just Euclidean)
- Smoothing and local linear regression can improve performance (at the cost of higher variance)
- With a lot of data, “local methods” have strong, simple theoretical guarantees
- Without a lot of data, neighborhoods aren’t “local” and methods suffer (curse of dimensionality)

→ more likely in large feature space