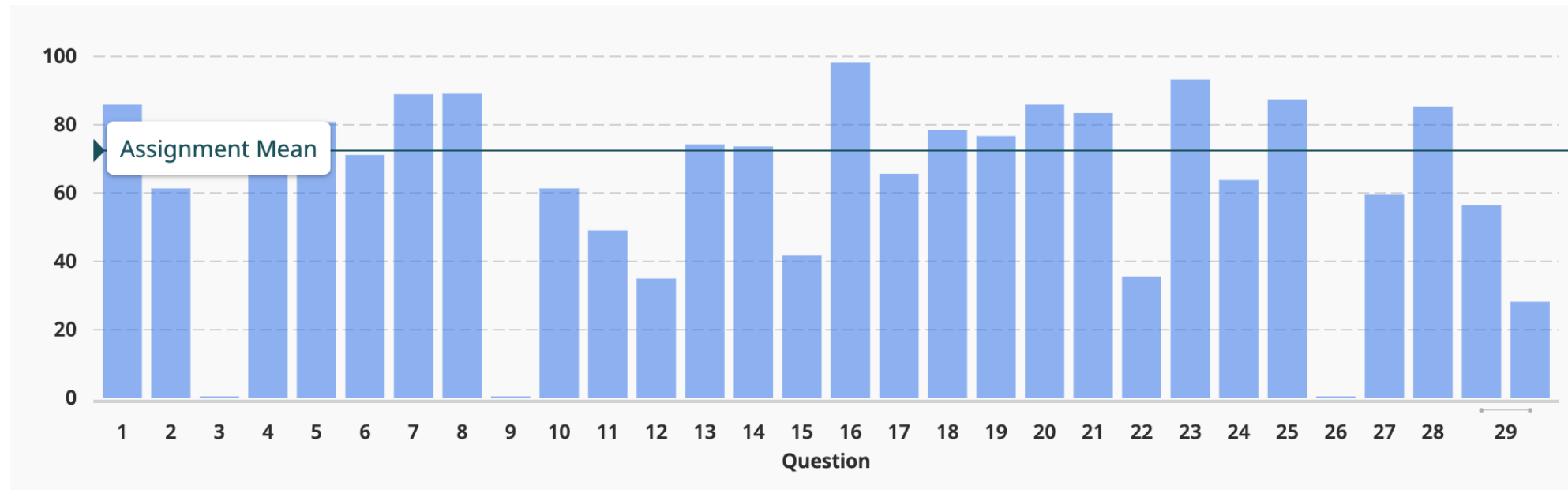


Midterm grades adjusted

- 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



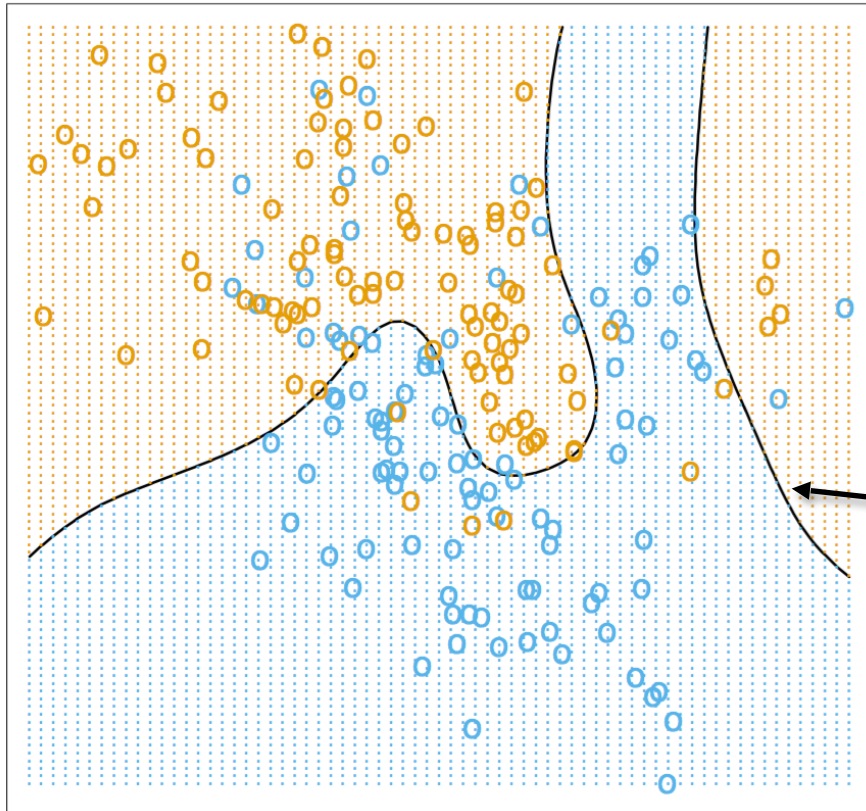
Parametric vs non-parametric

- A model is parametric if # parameters does not depend on # samples
- A model is non-parametric if # parameters increases with # samples
 - Does not mean absence of parameters!

This lecture: k nearest neighbors

- Assume we have a classification task
- To classify a new point x :
 - Find its k nearest neighbors in the training data
 - 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

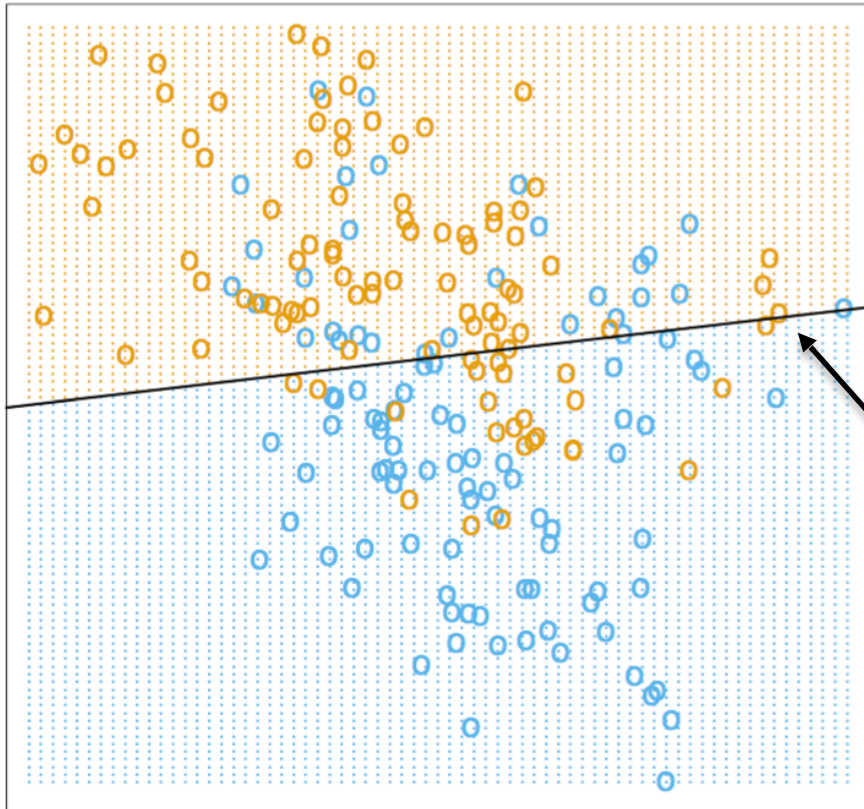
Optimal Bayes classifier:

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

▨ Predicted label: +1

▨ Predicted label: -1

Linear decision boundary



Training data:

○ True label: +1

○ True label: -1

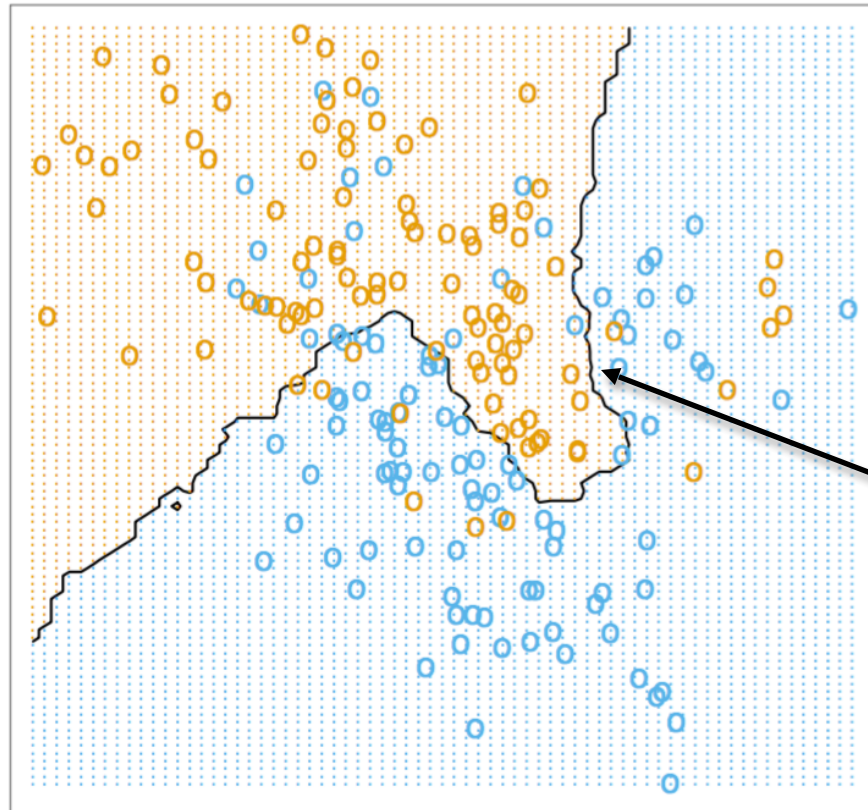
Learned linear decision boundary:

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

▨ Predicted label: +1

▨ Predicted label: -1

$k = 15$ nearest neighbors boundary



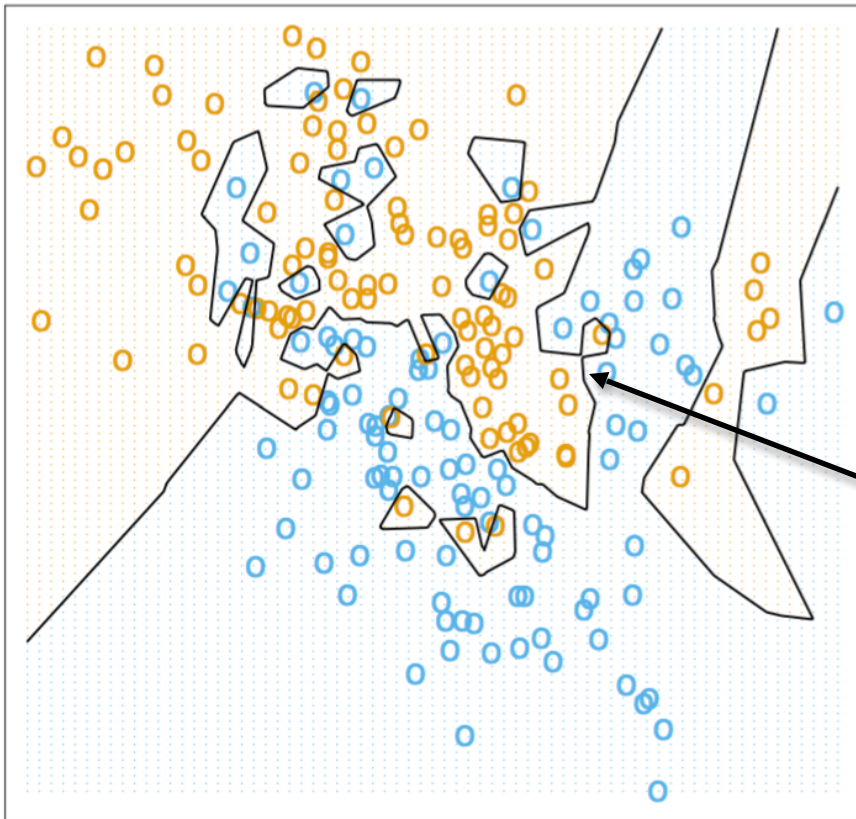
Training data:

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

15 nearest neighbors
decision boundary (majority
vote)

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

$k = 1$ nearest neighbor boundary



Training data:

○ True label: +1

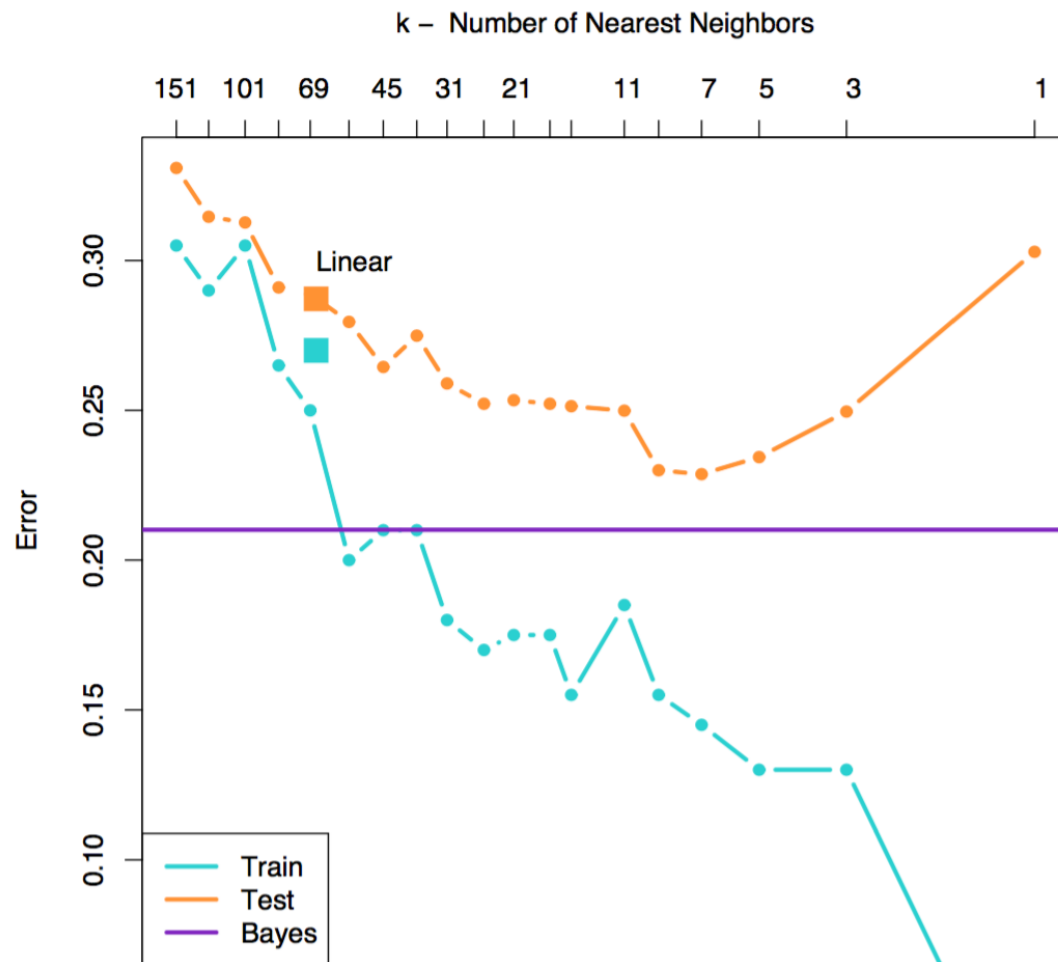
○ True label: -1

1 nearest neighbor decision boundary (majority vote)

■ Predicted label: +1

■ Predicted label: -1

k nearest neighbors error

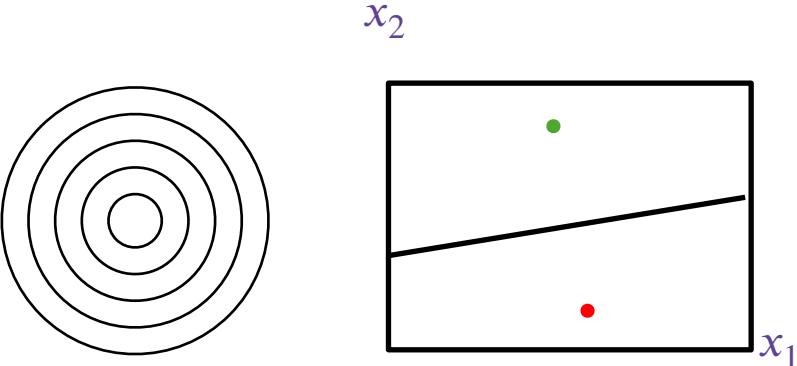


Parametric vs non-parametric

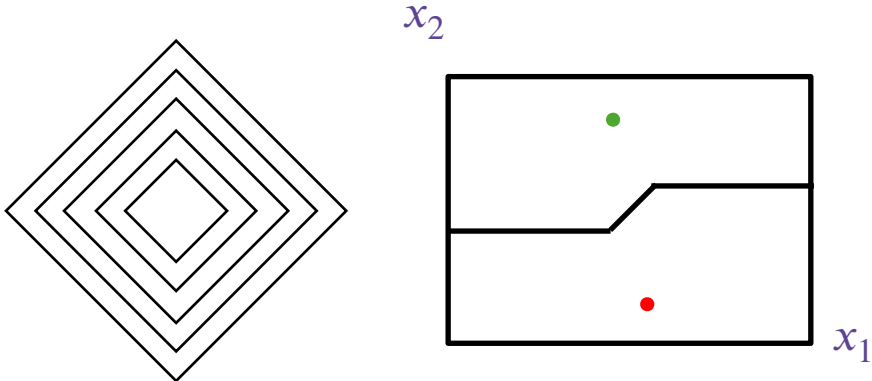
- A model is parametric if # parameters does not depend on # samples
- A model is non-parametric if # parameters increases with # samples

Notable distance metrics & level sets

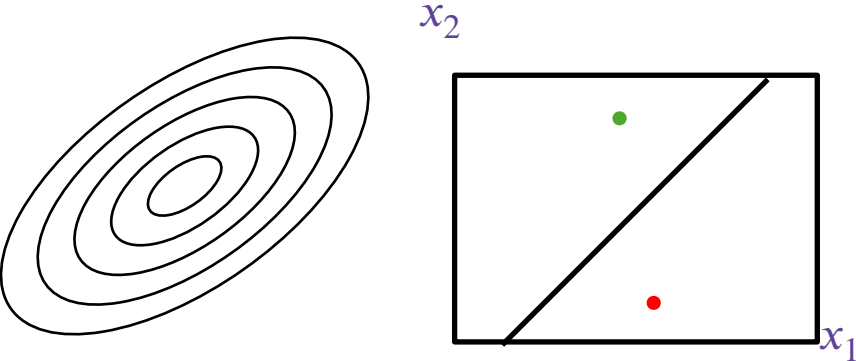
ℓ_2 norm (Euclidean)



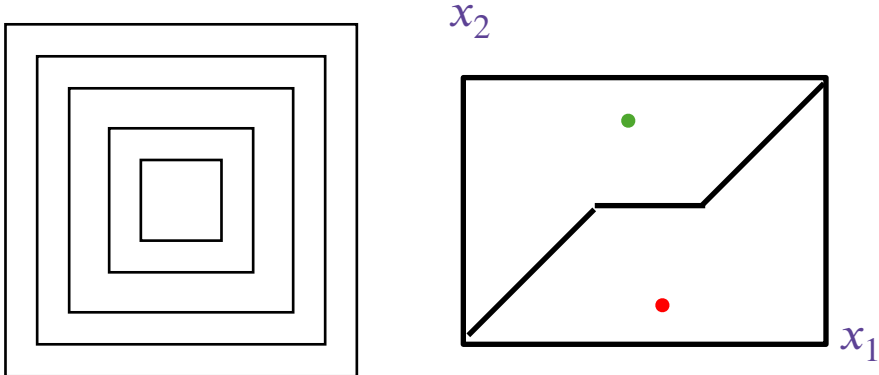
ℓ_1 norm (Manhattan, taxicab)



Mahalanobis norm

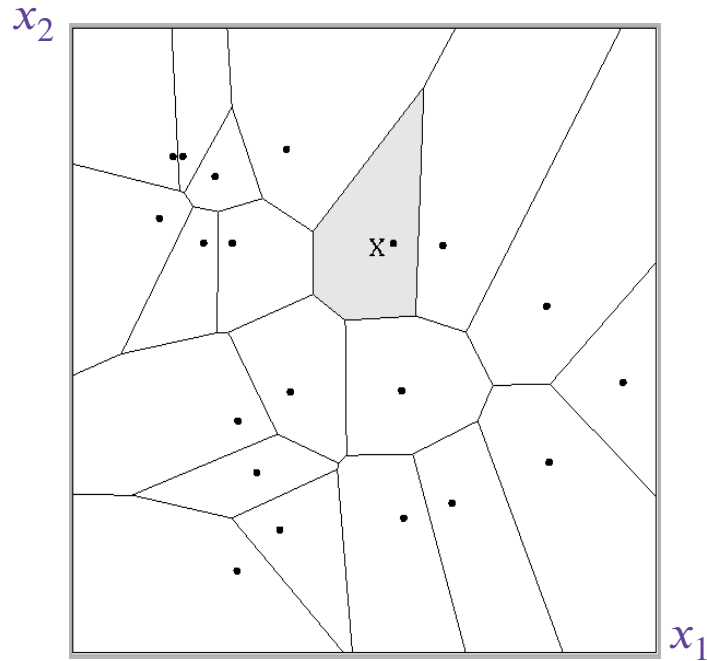


ℓ_∞ norm (max)

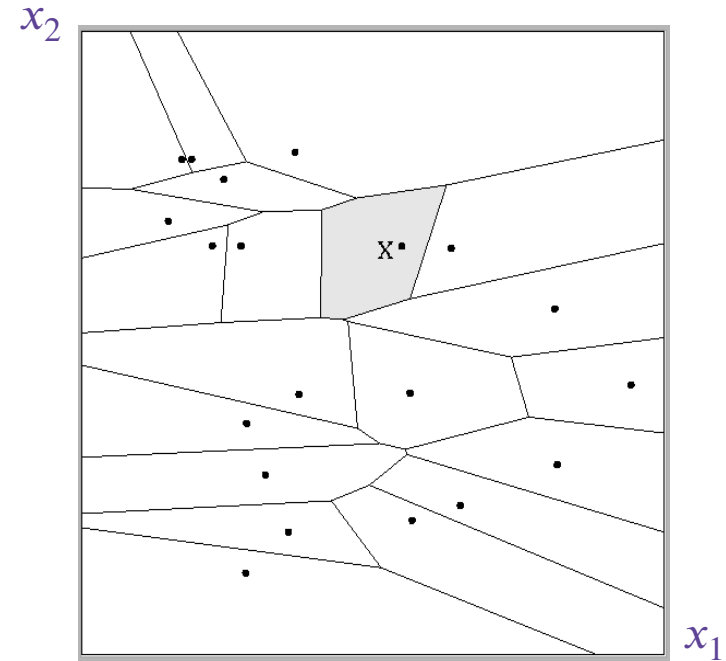


Example: distance metrics with $k = 1$ NN

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



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



Learned distance metrics

Training data



Dog



Cat

Test data

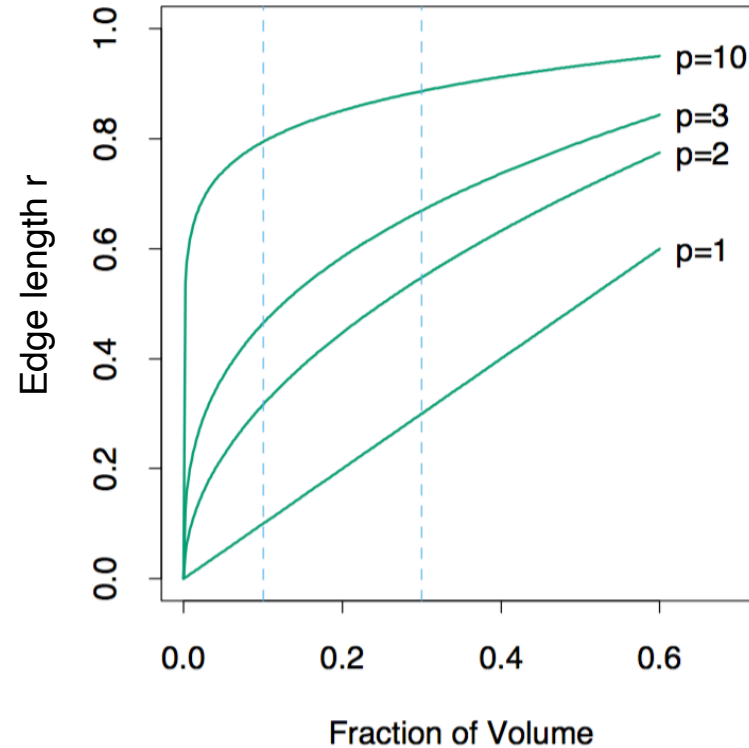
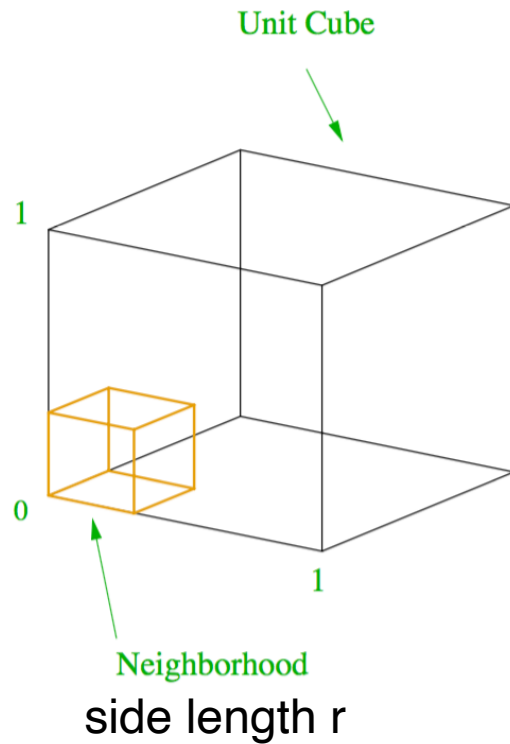


1-NN classification: Theoretical guarantees

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.

Curse of dimensionality, example 1

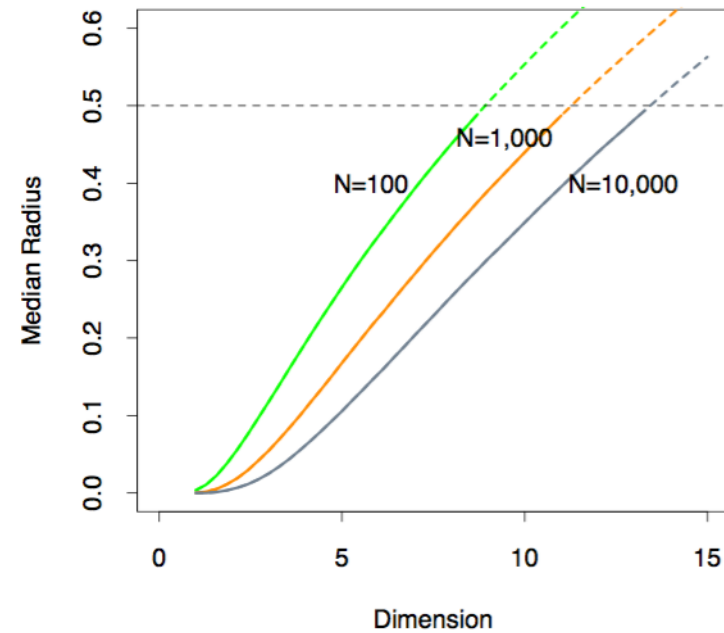
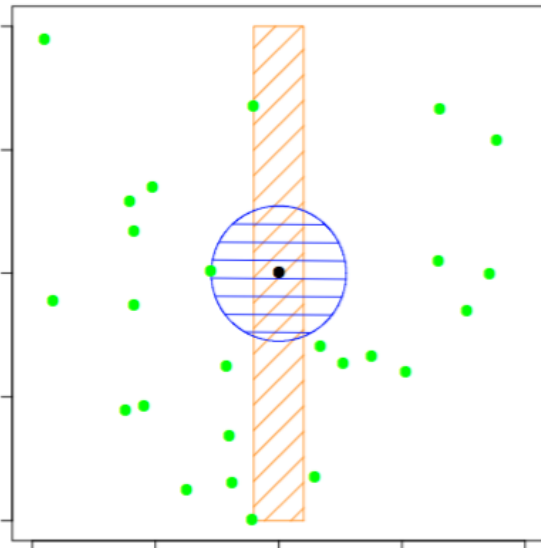


X is uniformly distributed over $[0, 1]^p$. What is $\mathbb{P}(X \in [0, r]^p)$?

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

Curse of dimensionality, example 2

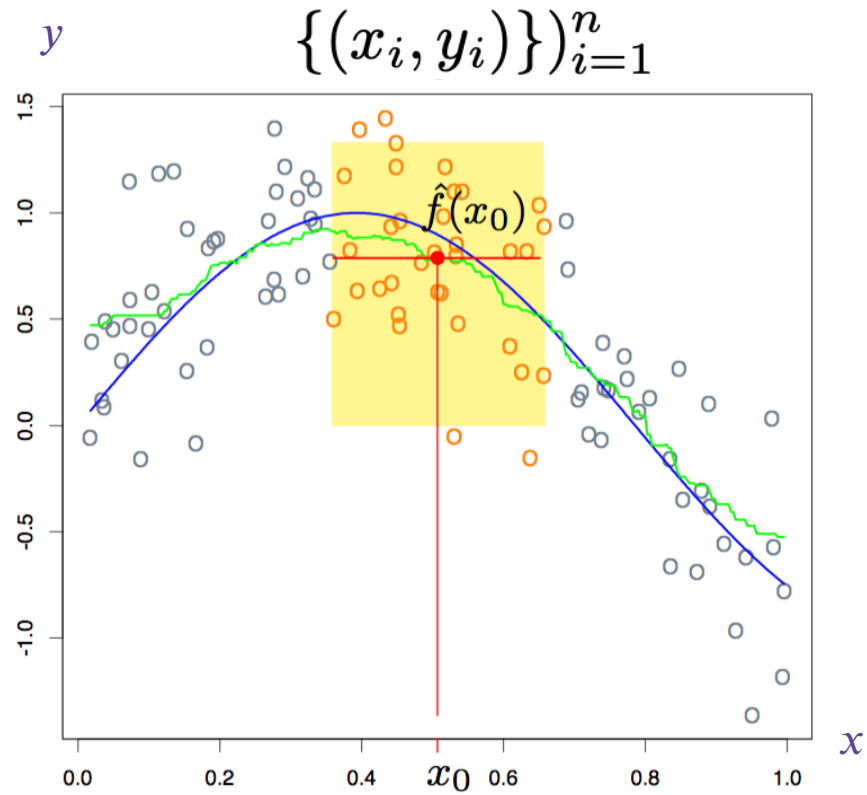
$\{X_i\}_{i=1}^n$ are uniformly distributed over $[-.5, .5]^p$.



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

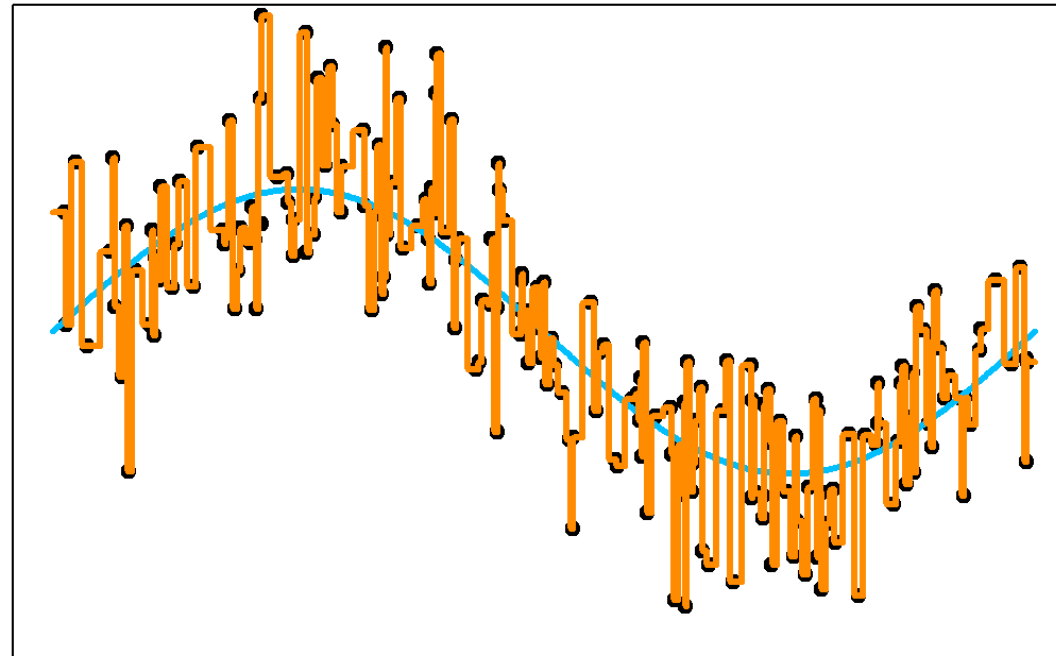
How many samples do we need so that a median Euclidean distance is within r ?

Nearest neighbor regression

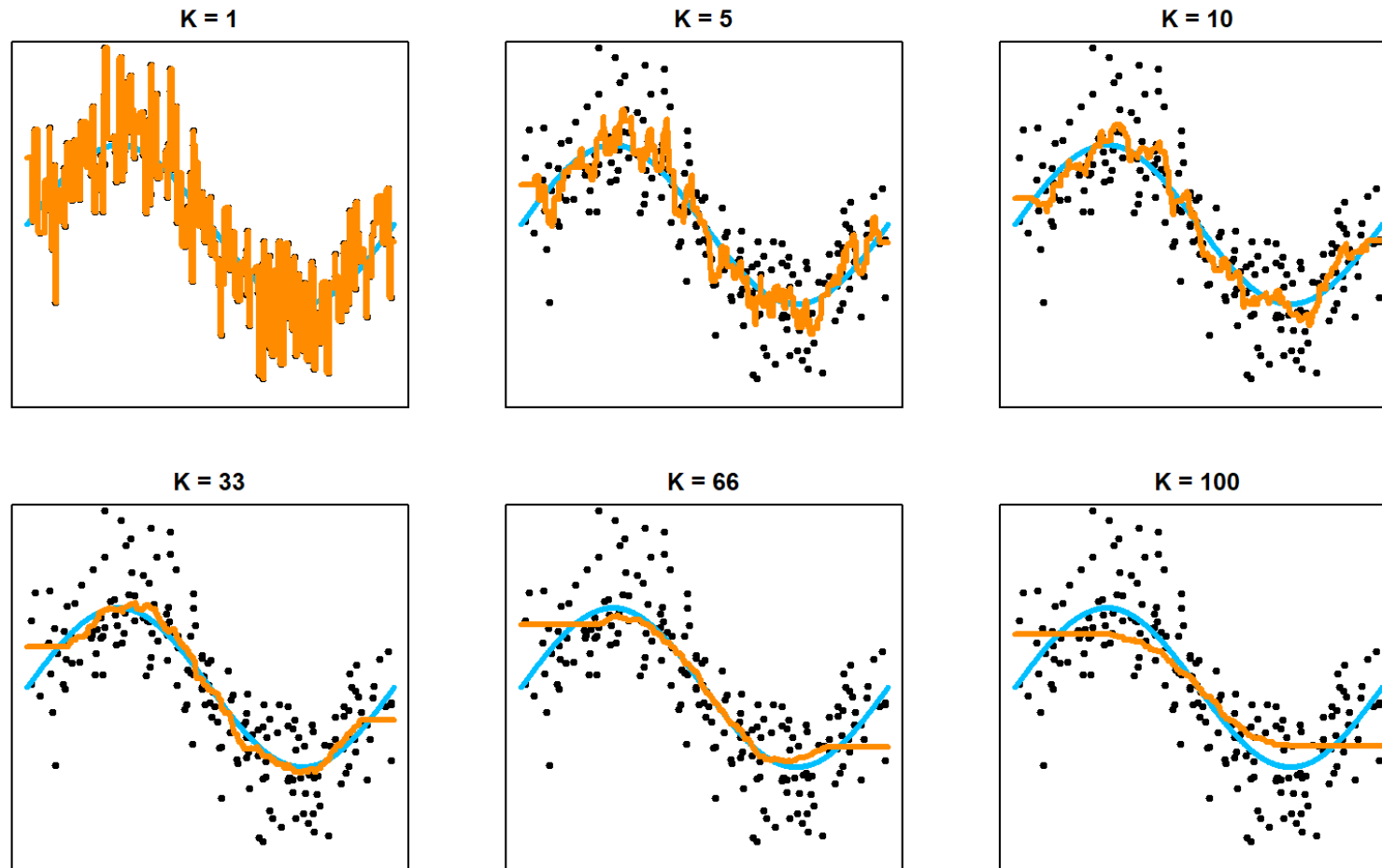


Overfitting

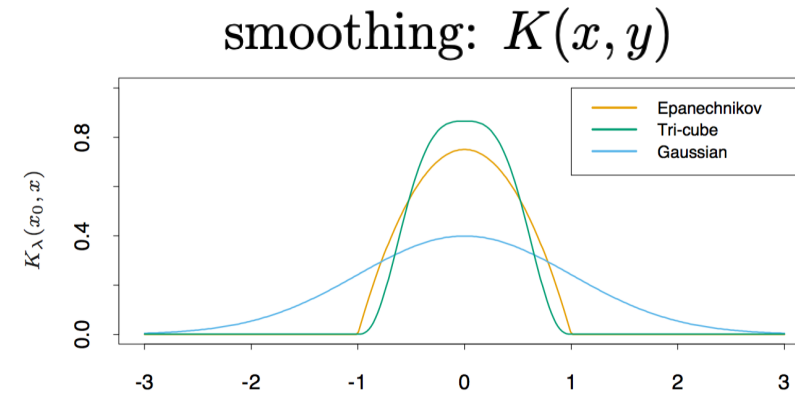
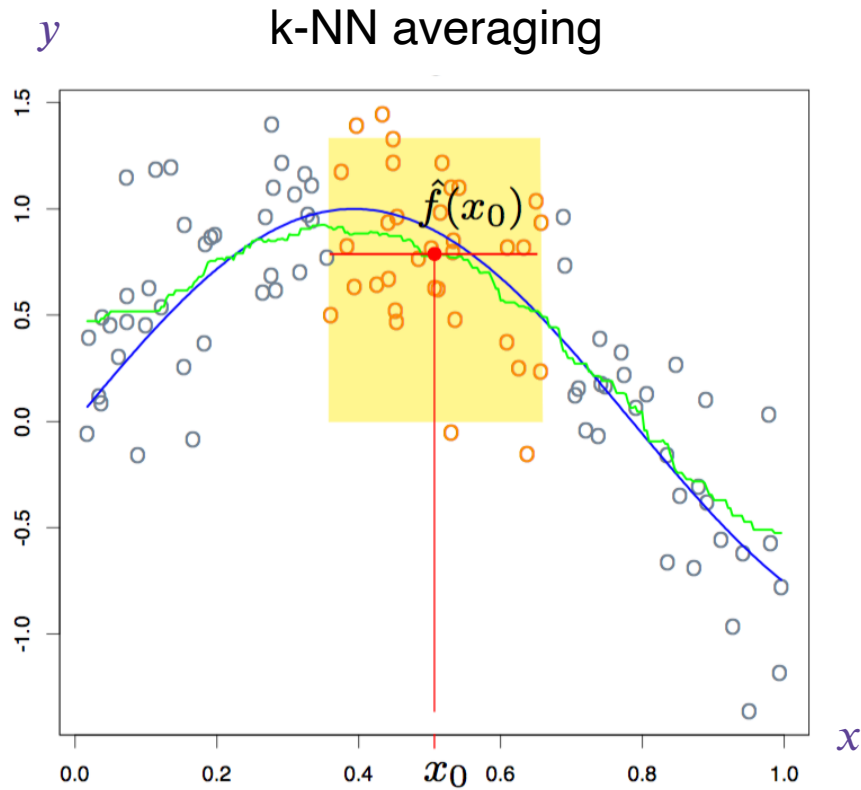
1-Nearest Neighbor Regression



Bias vs variance



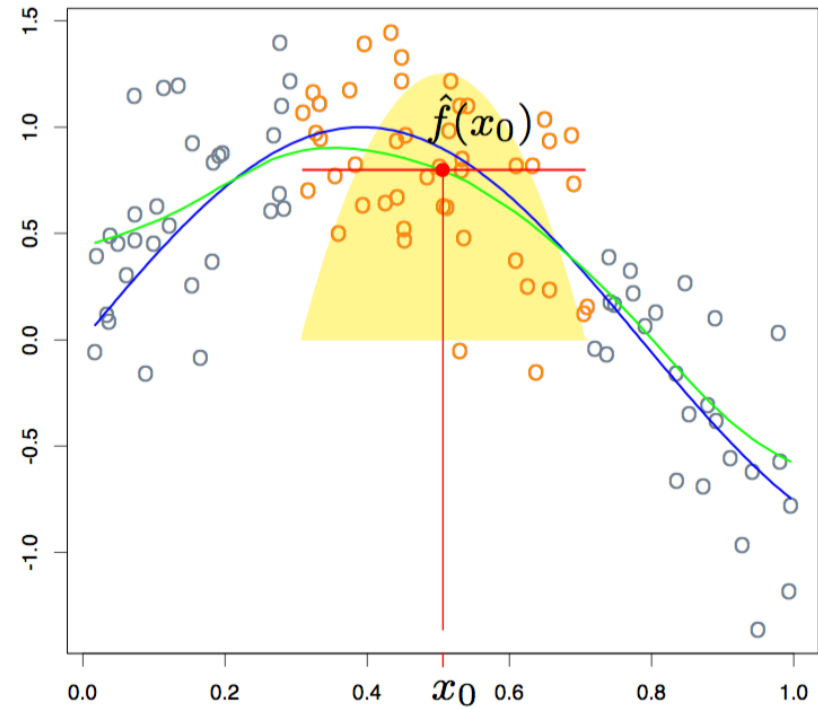
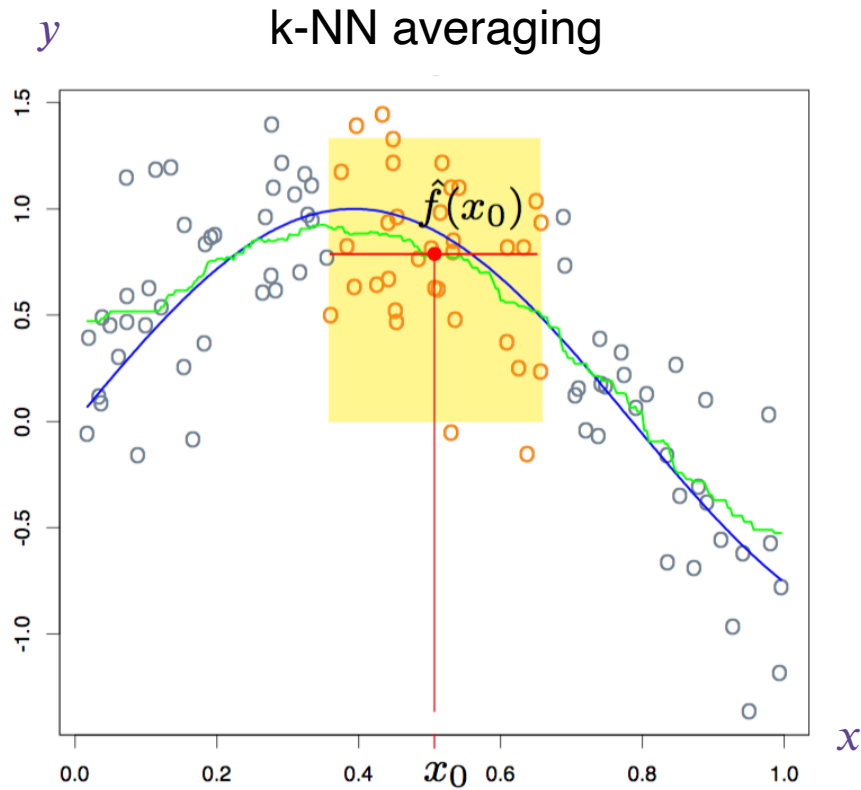
Smoothed nearest neighbor regression



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

Smoothed nearest neighbor regression

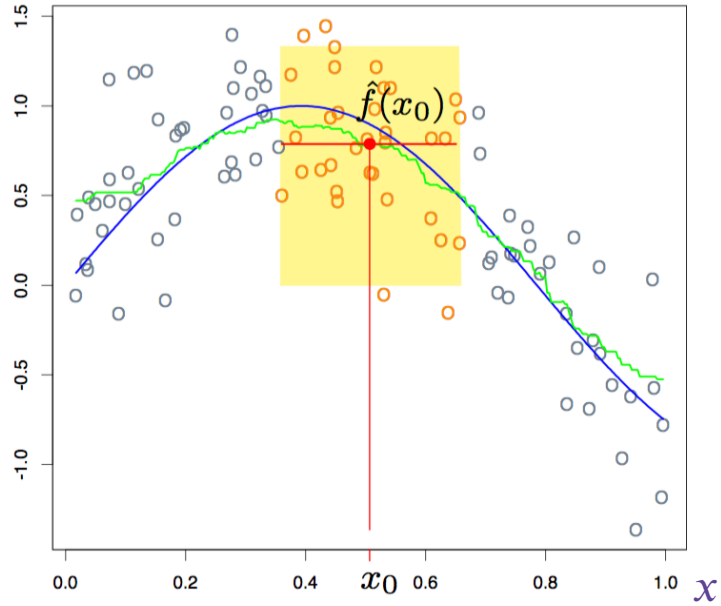
$$\hat{f}(x) = \frac{\sum_{i=1}^n K(x, x^{(i)}) y^{(i)}}{\sum_{i=1}^n K(x, x^{(i)})}$$



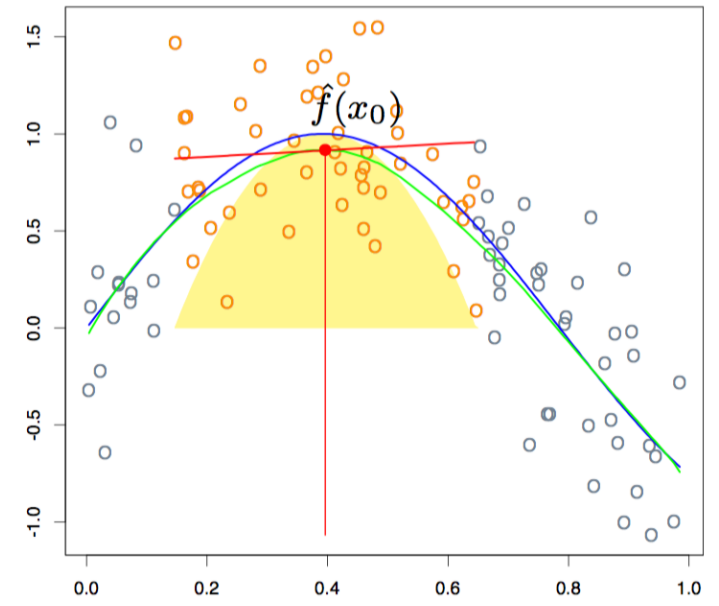
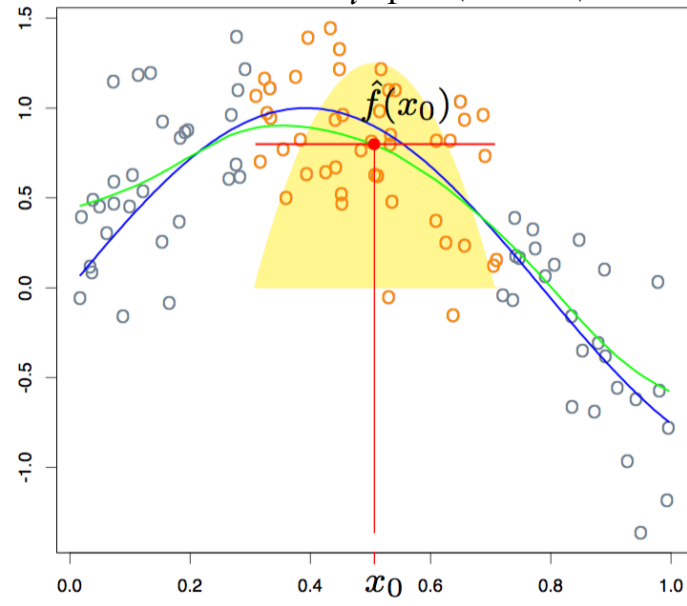
Locally linear regression

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:

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

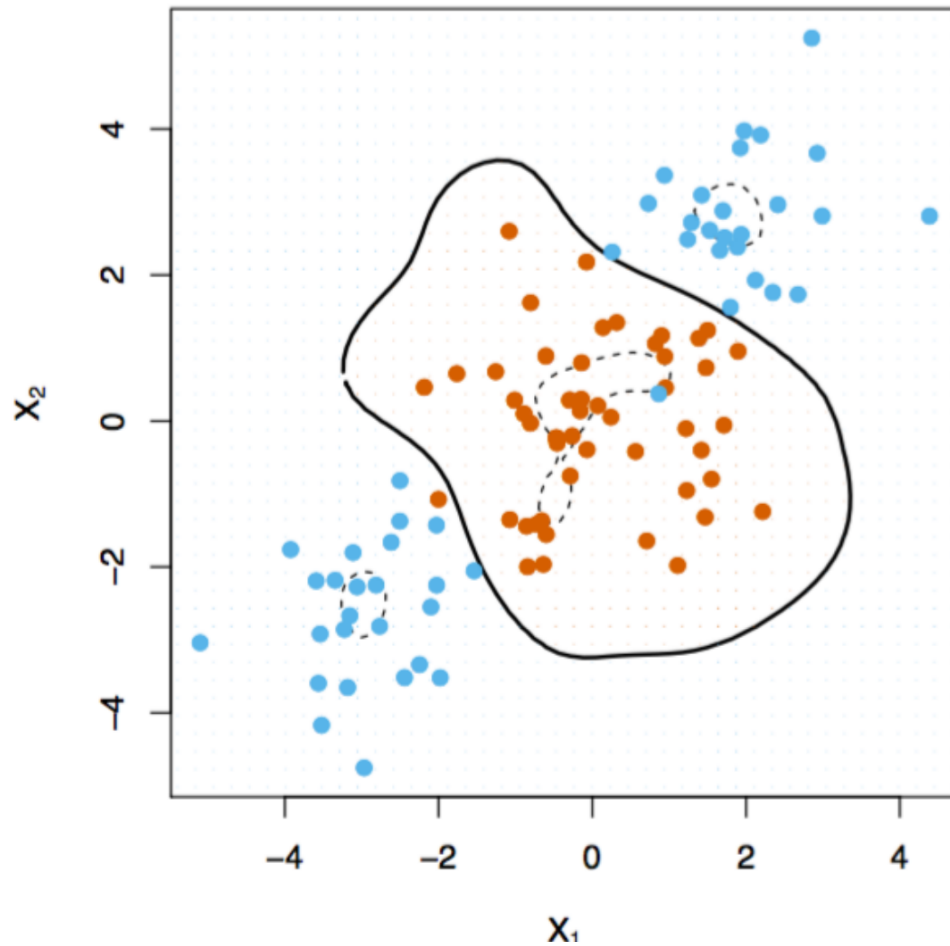
parameters goes up with # data

- Compare with (smoothed) nearest neighbors:

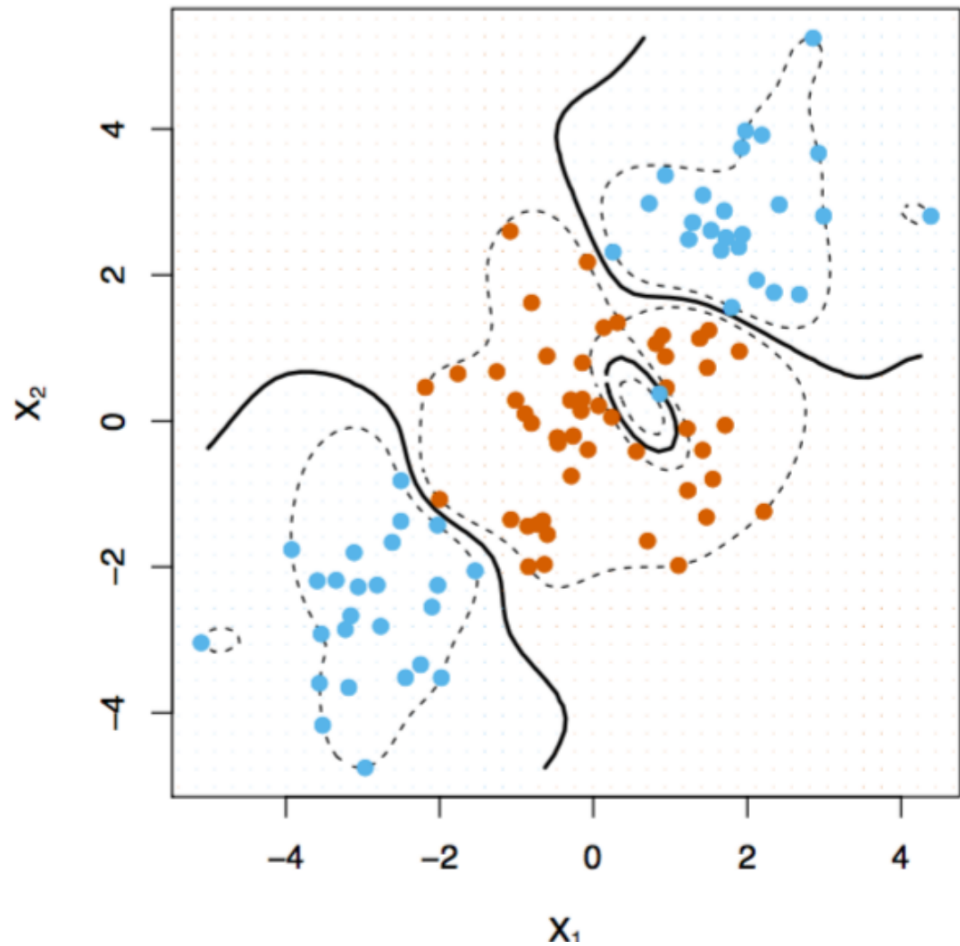
$$\hat{f}(x) = \frac{\sum_{i=1}^n K(x, x^{(i)}) y^{(i)}}{\sum_{i=1}^n K(x, x^{(i)})}$$

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

Bandwidth σ is large enough

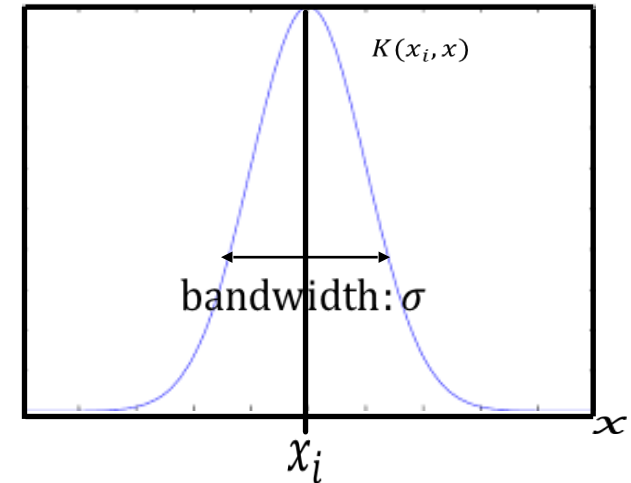
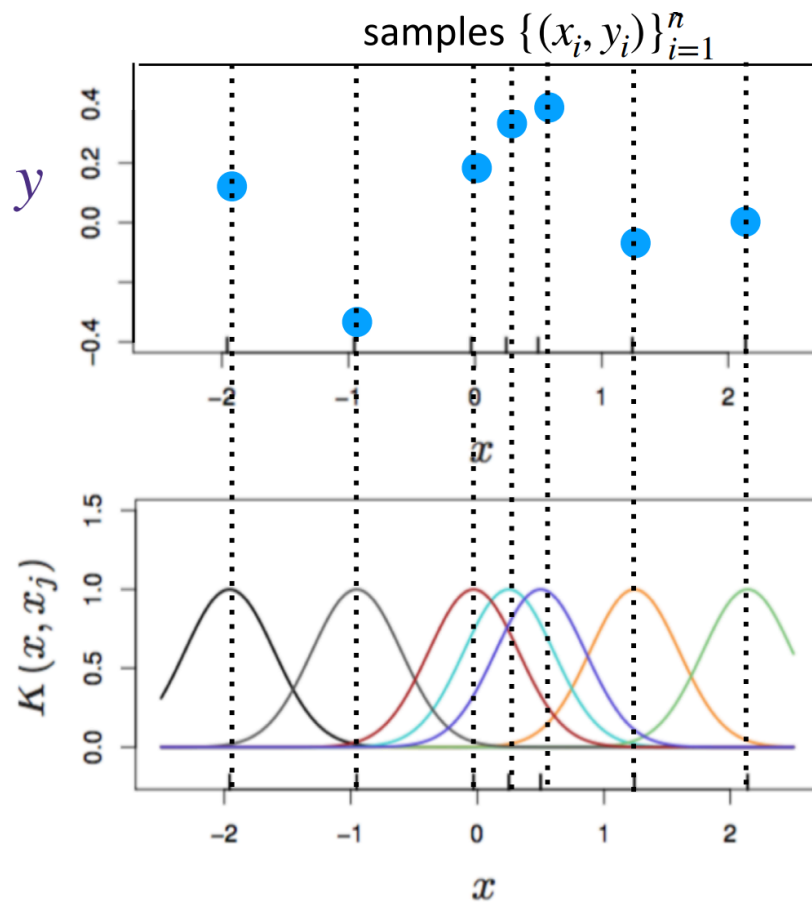


Bandwidth σ is small



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?

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)