

Today

- Classification case

Regression & classification

- Regression:
 - numerical response variable
 - predict a numerical value given x
 - e.g. number of species given latitude
- Classification:
 - categorical response variable
 - predict the category given x
 - e.g. is it a bird, deer, tree, or mountain lion?
 - e.g. is it dead or alive?; present or absent?

Classification

As before: out-of-sample accuracy

One common, simple measure is the error rate. If we have a *test* dataset of $i = 1 \dots n$ observations, the out-of-sample error rate is:

$$\frac{1}{n} \sum_i^n I(y_i \neq \hat{y}_i) = \text{mean}(I(y_i \neq \hat{y}_i)) \quad = \text{proportion incorrect}$$

\hat{y}_i is the predicted category for test case i .

$I()$ is an indicator function that equals 1 if the prediction is *incorrect* (i.e. if $y_i \neq \hat{y}_i$) and 0 if the prediction is correct.

Theory: best we can do

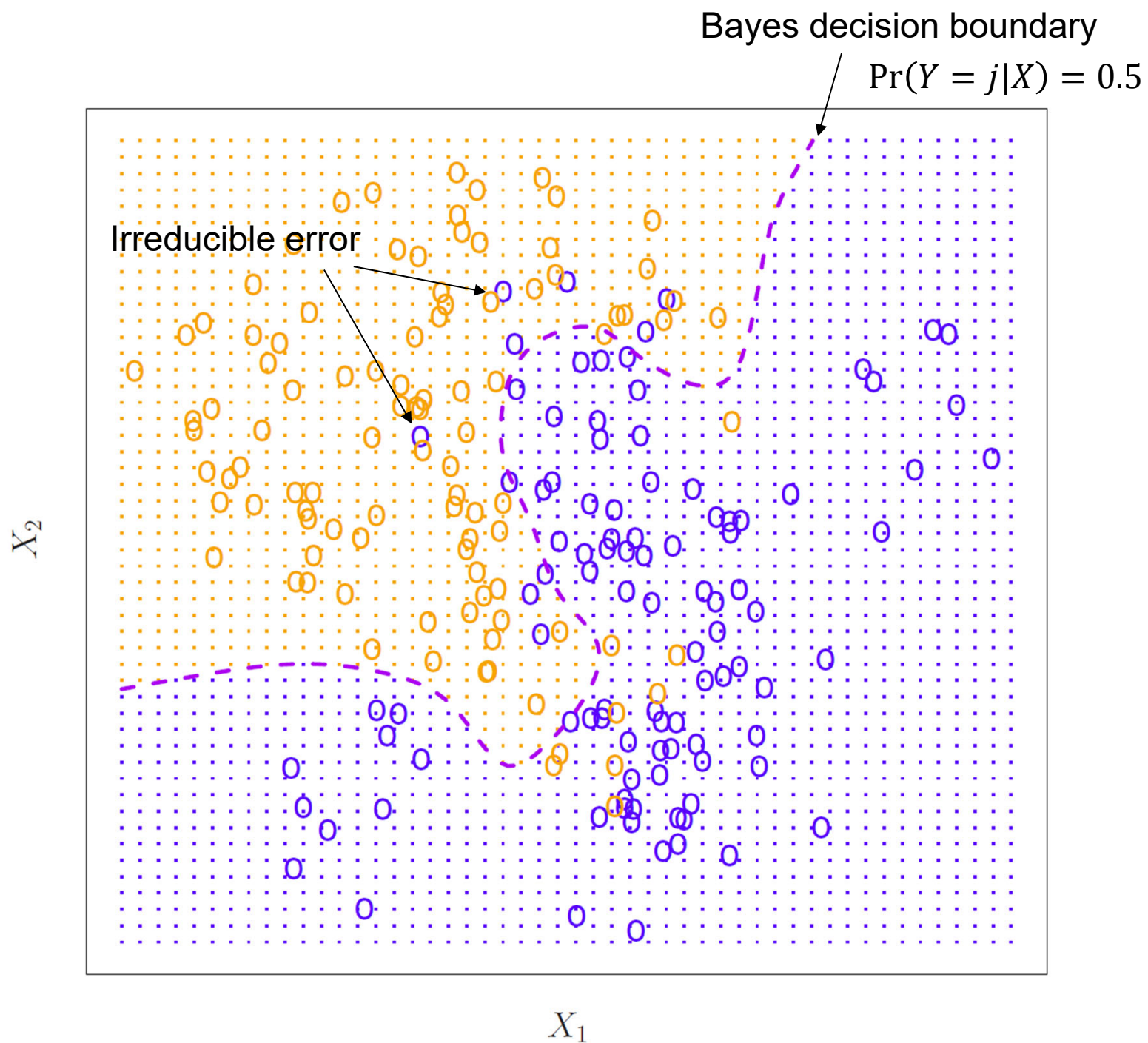
(assuming we know the true data generating process, e.g. in the case of data simulated from a model)

Bayes classifier

We should predict that our new observation belongs to the category with the highest probability. That is, choose category j for which

$$\Pr(Y = j | X = x_{\text{new}})$$

is highest.



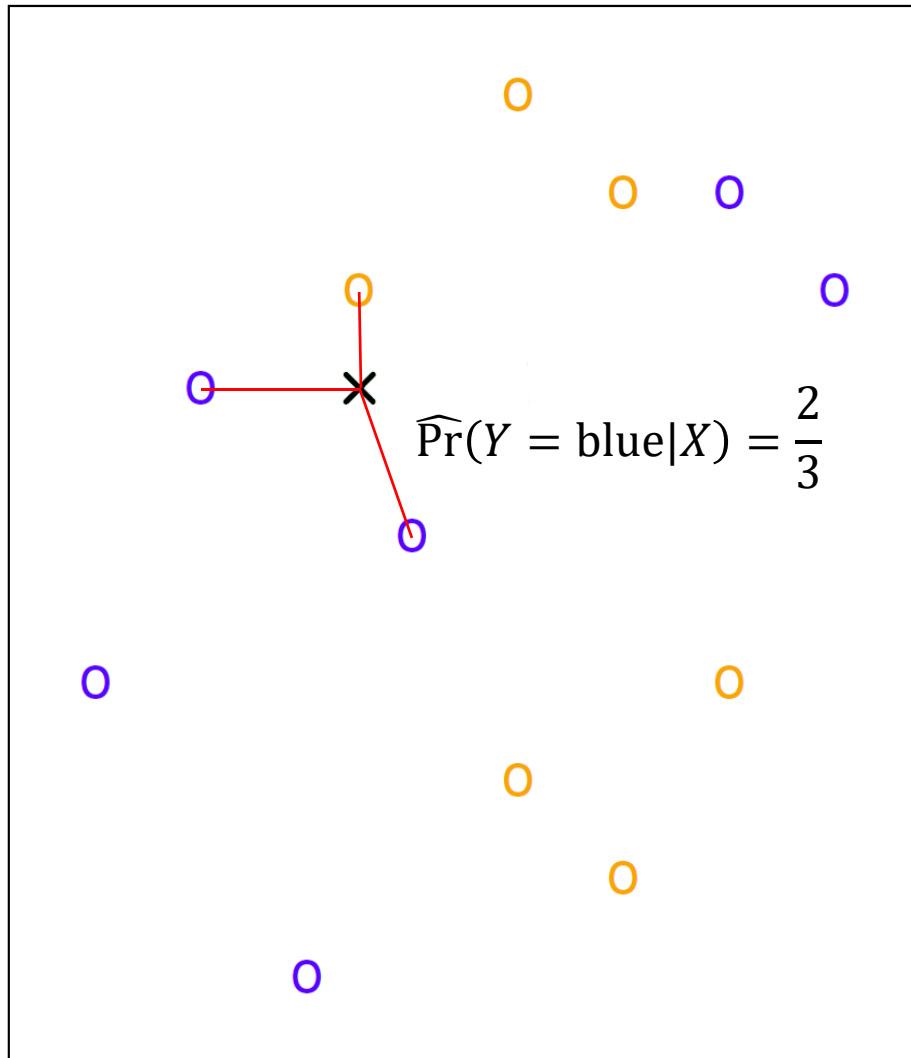
Theory: best we can do

Bayes error rate (irreducible error)

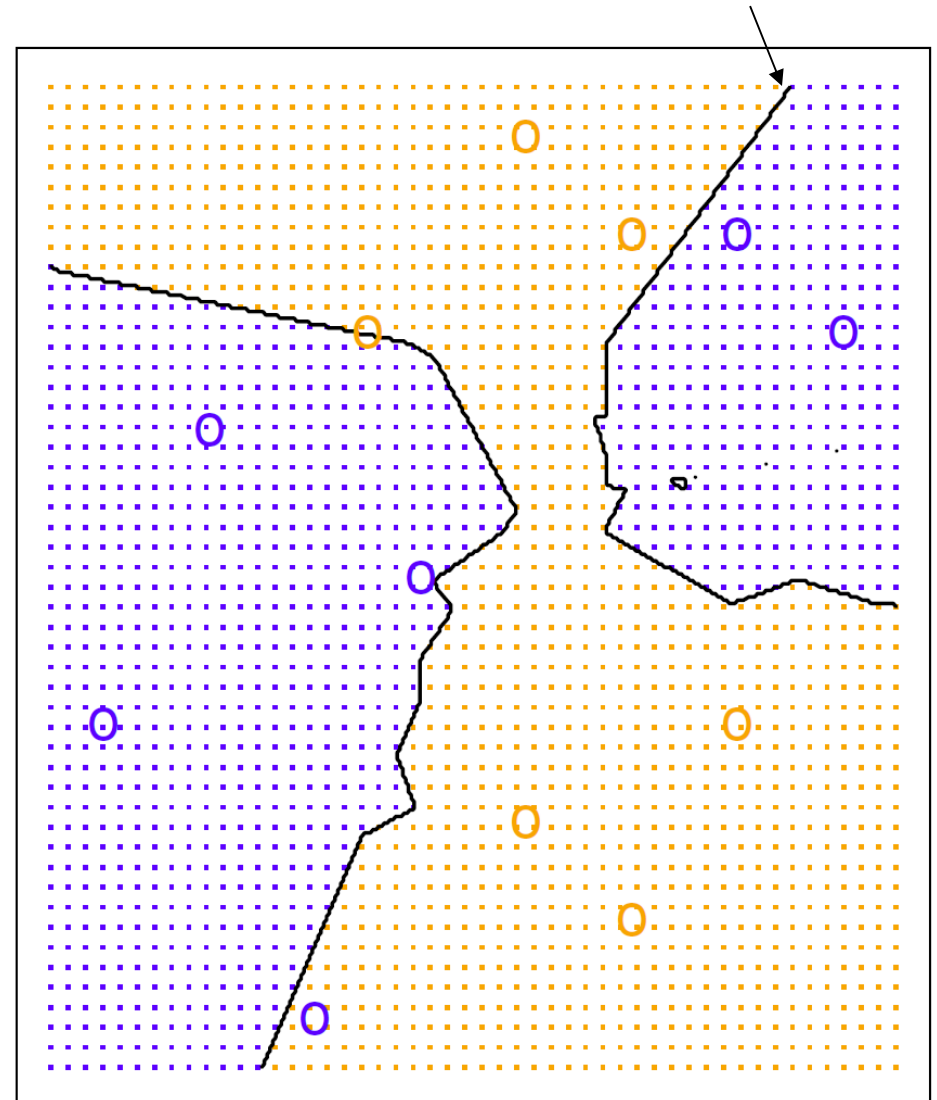
$$1 - E \left(\max_j \Pr(Y = j|X) \right)$$

E = expected value = mean over all values of X .

KNN
 $k = 3$



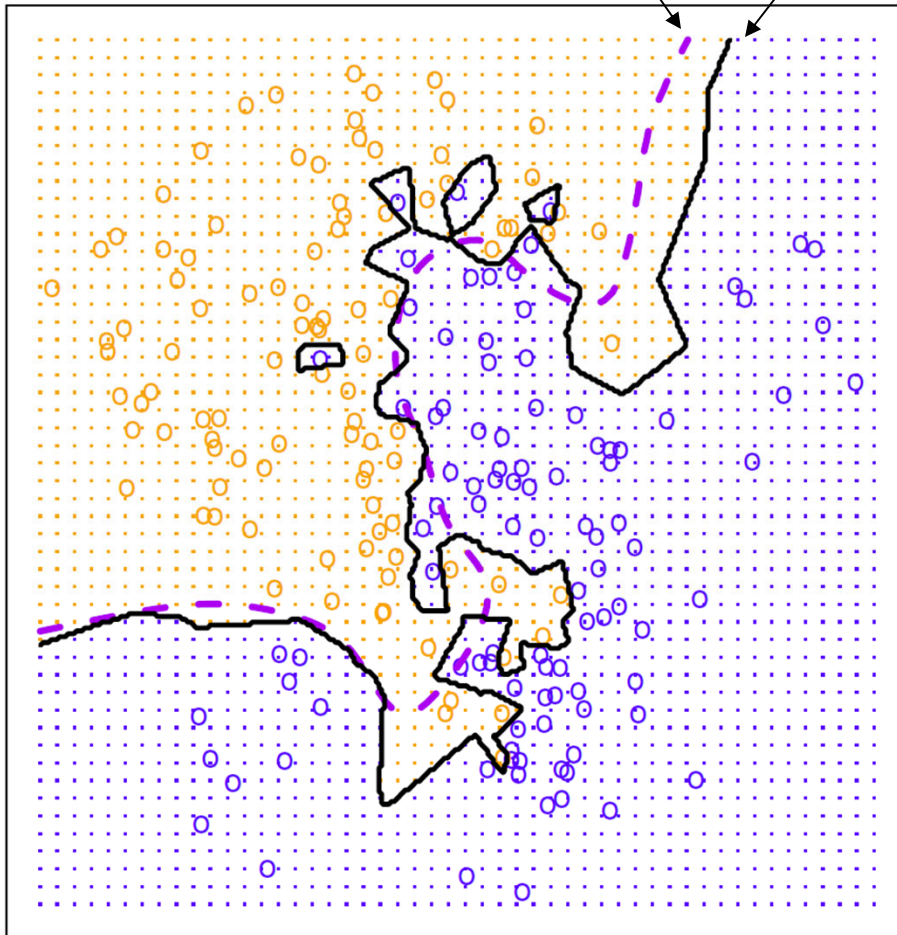
KNN decision boundary
 $\widehat{\Pr}(Y = j|X) = 0.5$



KNN
 $k = 1$

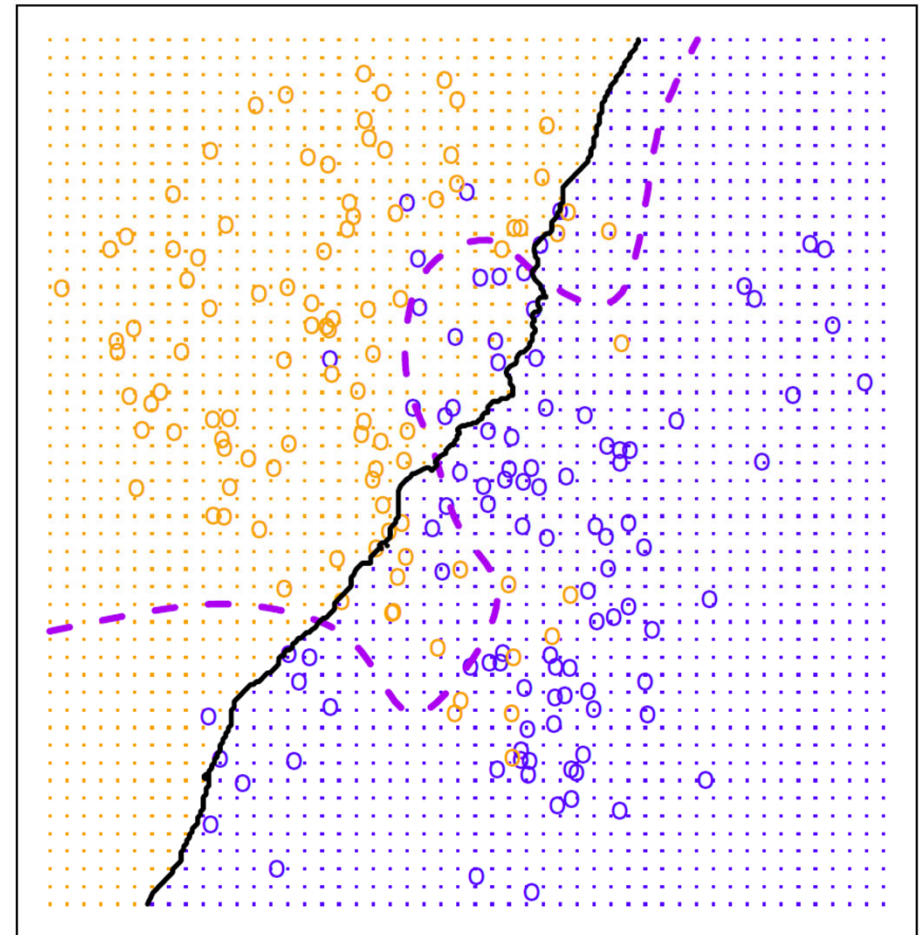
Bayes decision boundary

KNN decision boundary



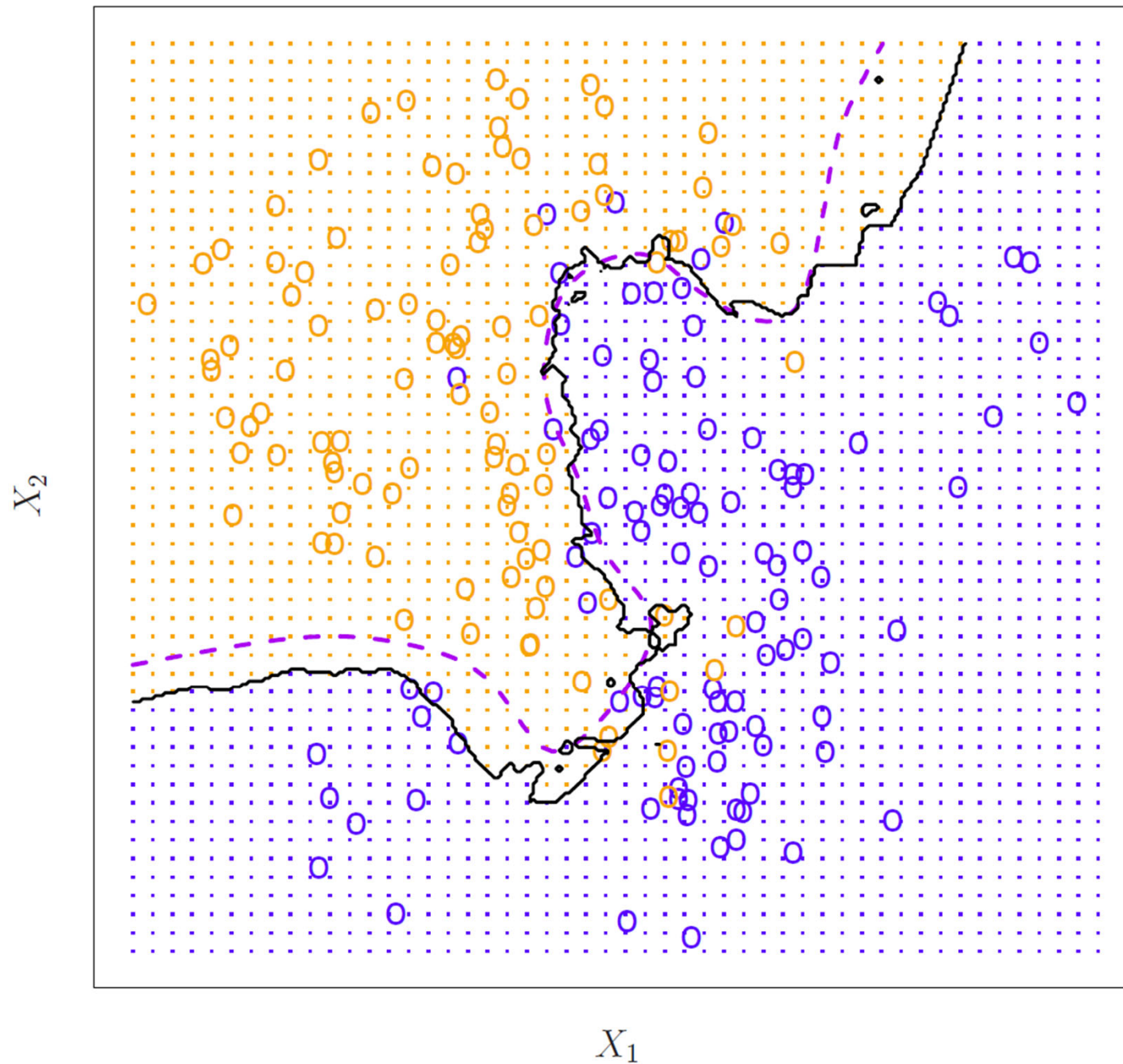
Overfit
high variance

KNN
 $k = 100$



Underfit
high bias

KNN
 $k = 10$



Mean across 250 k-fold CV runs

