

Lesson 6 - 12/09/24

Classification Problem:

- From X we predict one of K classes.
- E.g.

A person with certain symptoms could have:

- Stroke $\rightarrow C_1$
- Diabetic ketoacidosis $\rightarrow C_2$
- Epileptic Seizure $\rightarrow C_3$

Which of the three does the individual have?

Bayes in Classification

- Look at X to predict C_k .
- As with Logistic Regression, useful to consider probs:
 $p(C_k|X)$.
By knowing this, with new datapoint, we can predict C_k
- We can estimate $p(C_k|X)$ from training data using Baye's Theorem:

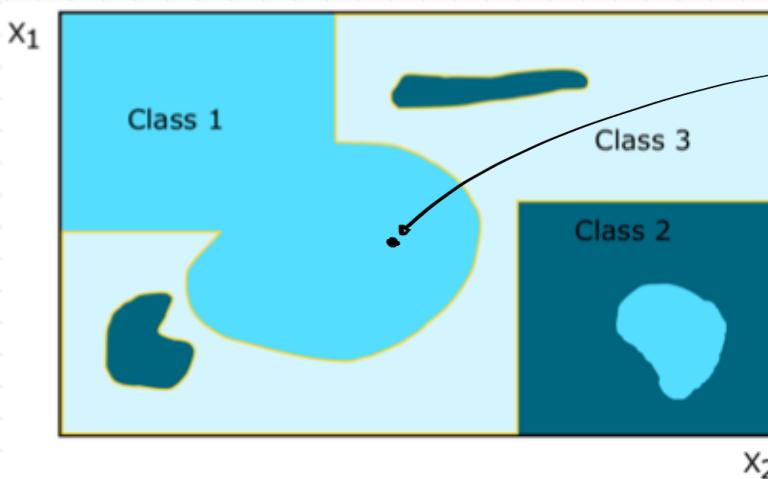
$$p(C_k|X) = \frac{p(X|C_k) p(C_k)}{p(X)}$$

- Let's use $p(C_k|X)$ to classify.
Goal: as few misclassification as possible.

WE NEED DECISION RULES!!!

- When do we classify a point as C_k ?
- Any decision rule divides feature space into design regions.
- Decision regions are separated by decision boundaries!

D.R.



IF A DATAPPOINT FALLS HERE, IT WILL BE PREDICTED AS C_1

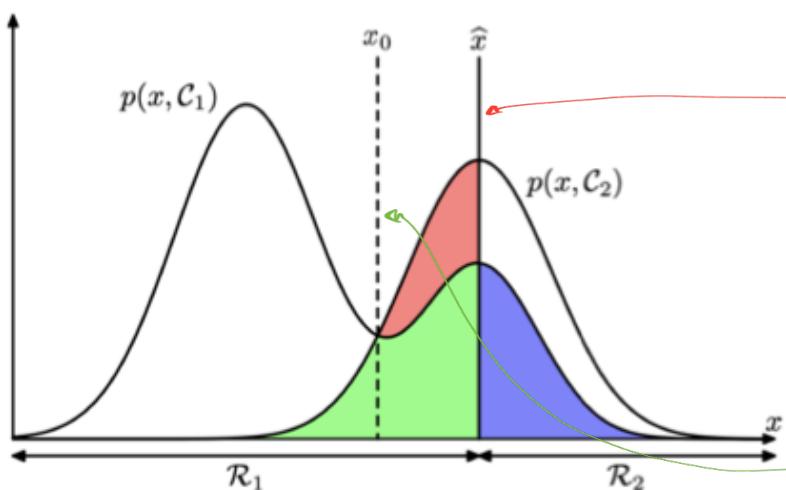
Aim to make as few misclassifications as possible. Minimize:

$$\begin{aligned} p(\text{mistake}) &= p(X \in \mathcal{R}_1, C_2) + p(X \in \mathcal{R}_2, C_1) \\ &= \int_{\mathcal{R}_1} p(X, C_2) dX + \int_{\mathcal{R}_2} p(X, C_1) dX \end{aligned}$$

From the product rule, $p(X, C_k) = p(X)p(C_k|X)$, so we get,

$$p(\text{mistake}) = \int_{\mathcal{R}_1} p(X)p(C_2|X) dX + \int_{\mathcal{R}_2} p(X)p(C_1|X) dX$$

$p(x)$ contributes to both integrals so mistake is limited the most by the assignment on the class with the highest $p(C_k|X)$ at x .



- Look at the Areas under the Functions!
- To the left classify as C_1 , to the right C_2

- Here we correctly classify C_2 but missclassify C_1
- Correctly classifying C_1 but missclassifying C_2
- Missclassifying everything

That's why we move the boundary here
MINIMIZE INTEGRALS!!!!

- What we've just developed is the Bayes classifier.
- With Bayes classifier we classify to K with highest posterior probs:

$$d(x) = \arg \max P(Y=y|x)$$

• NOT ALL MISTAKES ARE EQUAL THOUGH

- There are some that are worse than others!
- I.e. for cancer, a **False Positive** causes **stress** to patient
a **False Negative** may cause the **death** of a patient!
- We really need to limit **False Negatives**!

Loss & Loss Matrix!

- We want to **punish** the algorithm **&** classification it makes
- We **assign** a loss $l_{kj} \geq 0$ depends on K , true class, and j , assigned class.
 - l_{kj} is an entry in a loss matrix
- The **bigger** the loss for a given classification, the **bigger** the punishment!

	cancer	normal	
cancer	0	1000	NO PUNISHMENT
normal	1	0	BIG PUNISHMENT!

THE CONCEPT OF LOSS LET US EXPRESS OUR GOAL:

- WE WANNA MINIMIZE THE EXPECTED LOSS.

$$\mathbb{E}[L] = \sum_k \sum_j \int_{\mathcal{R}_j} L_{kj} p(X, C_k) dX$$

Again, we can rewrite this to get

$$\mathbb{E}[L] = \sum_k \sum_j \int_{\mathcal{R}_j} L_{kj} p(X) p(C_k | X) dX.$$

So we minimize loss by assigning a new X to the class j that minimizes

$$\sum_k L_{kj} p(C_k | X).$$

MINIMIZING POSTERIOR EXPECTED LOSS IS ENOUGH

- IF $d(x)$ MINIMIZES POSTERIOR EXPECTED LOSS $\forall x$
THEN $d(x)$ ALSO MINIMIZES THE EXPECTED LOSS

- WE OFTEN USE 0-1 LOSS:

$$d(j, k) = \begin{cases} 1, & j \neq k \\ 0, & j = k \end{cases}$$

In regression, we can also use loss, e.g. *Squared error loss*

$$L(y, d(x)) = (y - d(x))^2.$$

Absolute error loss

$$L(y, d(x)) = |y - d(x)|.$$

Remember: the loss function is usually non-negative.

THE EXPECTED LOSS IS A THEORETICAL QUANTITY:

ALSO NAMED → TEST ERROR, GENERALISATION ERROR, RISK, PREDICTOR ERROR.

- WE CAN ESTIMATE EXPECTED LOSS BY EMPIRICAL RISK:

$$\frac{1}{n} \sum_{i=1}^n L(Y_i, d(X_i)) \quad \text{WHERE } n \text{ IS THE # OF OBSERVATIONS}$$

TEST & TRAINING ERROR

- TRAINING ERROR IS THE EMPIRICAL RISK COMPUTED FROM THE TRAINING SET

- GENERALLY A BAD ESTIMATOR

- THE TEST ERROR DENOTES BOTH:

- THE TRUE EXPECTED LOSS

- THE ESTIMATE = EMPIRICAL RISK FROM TEST DATA

- TEST ERROR COMPUTED BY CROSS-VALIDATION IS ALSO ESTIMATE OF EXPECTED LOSS.

BAYES CLASSIFIER

IN THE LANGUAGE OF LOSS:

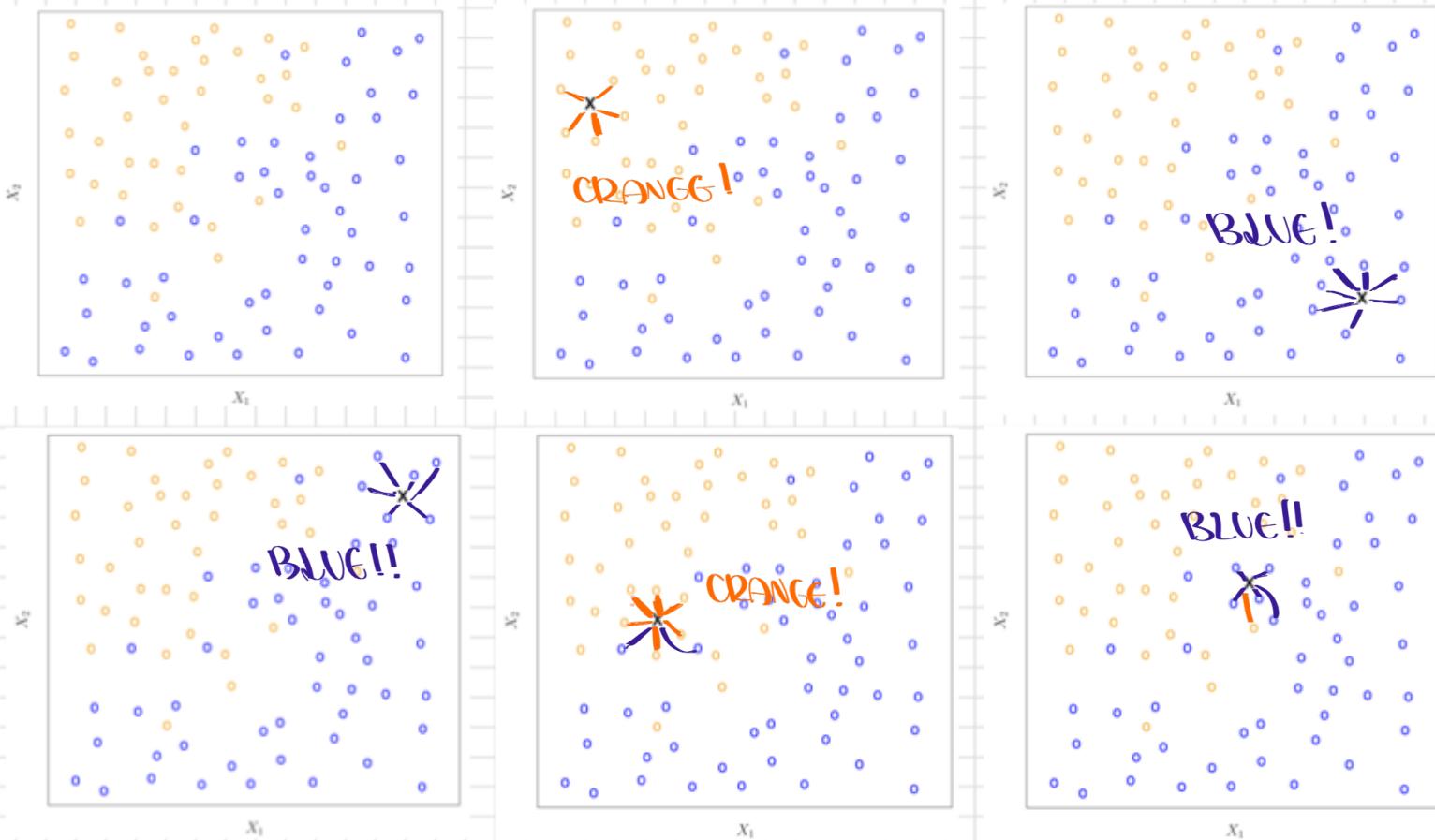
- THE BAYES CLASSIFIER MINIMISES THE EXPECTED LOSS UNDER THE SPECIFIC CHOICE 0-1 LOSS.

- The associated error, Bayes error rate, is a theoretical lower bound.
- Remember: with a Bayes classifier, we classify to the class K with the highest posterior probability.

$$d(x) = \arg \max P(y=j|x)$$

K-NEAREST NEIGHBOUR

- What if we see x 's that we've never observed? How do we classify?



- By looking at the picture, our first move is to look at its neighbours if then say that x will be part of that class.
- KNN (K-Nearest Neighbour) approximates the posterior class probs

Classify point x_0

- ① Find K points in the training data closest to $x_0 \rightarrow$ set N_0
- ② Estimate posterior probs for class j as fraction of N_0 from j :

$$P(Y=j|X=x_0) = \frac{1}{K} \sum_{i \in N_0} I(y_i=j).$$

- ③ Choose class with highest posterior probs.

- Choosing a right K matters, but obviously there is not a method to decide this!



KNN: K=1



KNN: K=100

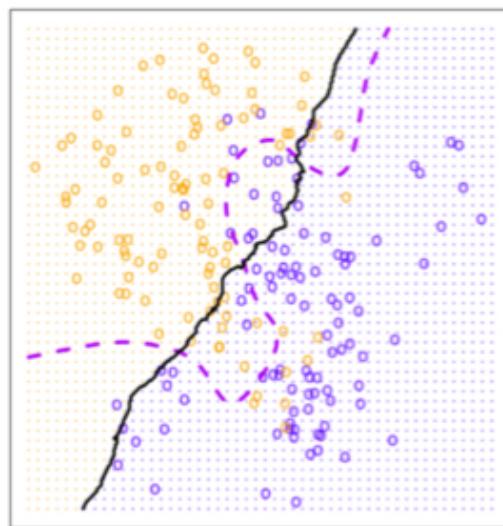


FIGURE 2.16. A comparison of the KNN decision boundaries (solid black curves) obtained using $K = 1$ and $K = 100$ on the data from Figure 2.13. With $K = 1$, the decision boundary is overly flexible, while with $K = 100$ it is not sufficiently flexible. The Bayes decision boundary is shown as a purple dashed line.

Summary for KNN:

- Resulting Decision Rule is simple:
 - KNN assigns a class according to a majority vote among K closest points.
- Gives extremely flexible boundaries
 - Recall image above.
- Good classifier with error rate close to Bayes error rate.
- Often does not work well in high dimension feature space.
 - many points, maybe big spread!
- K can be chosen by cross-validation