

Summary

This week, we focused on classification, where the goal is to predict qualitative outcomes based on input features. We discussed the difference between assigning a class label versus estimating the probability of belonging to a class, which can be more informative. We also examined the use of linear regression for classification and why logistic regression is often a better choice for producing valid probability estimates.

Concepts

Classification function: A model that maps input features X to qualitative outcomes Y .

Qualitative variables: Outcomes that fall into categories without an inherent order.

Probability estimation: It's more informative to estimate the probability that an input belongs to a class rather than just assign a label.

Credit card default example: Used to demonstrate classification based on balance and income.

Linear regression for classification: Codes the response as 0 or 1 and uses regression to estimate outcomes. It can predict values outside the $[0,1]$ range, making it unreliable for probability tasks.

Logistic regression: A more appropriate method for classification that outputs values between 0 and 1, modeling the probability directly.

Uncertainties

One area I'm still unsure about is when linear regression might be considered "good enough" for classification, despite its flaws. I'd also like to better understand how logistic regression is extended to handle more than two classes, as the lecture primarily focused on binary outcomes. Finally, I'm curious about real-world situations where having poor probability estimates has caused problems, which would help me appreciate why accurate estimates matter so much in practice.