

Classification: $Y_i = k, k \in \{1, \dots, C\}$ (special case: $C=2$ two-class classification)

• deterministic classifier: one hard decision for each i : $k \in \{0, 1\}$ $k \in \{-1, 1\}$

$$\hat{Y}_i = f(X_i) \text{ with } \hat{Y}_i = Y_i^* \text{ truth}$$

collect outcomes in "confusion matrix"

$\hat{Y} \backslash Y^*$	-1	+1
-1	true negative	false negative
+1	false positive	true positive

should be diagonal

ideally all are 0

quality metrics: (1) false positive fraction = $\frac{\#FP}{N} \in [0, 1]$

false negative fraction = $\frac{\#FN}{N} \in [0, 1]$

(2) false positive rate = $\frac{\#FP}{\#FP + \#TN} \approx P(\hat{Y}=1 | Y^*=-1)$

false negative rate = $\frac{\#FN}{\#FN + \#TP} \approx P(\hat{Y}=-1 | Y^*=1)$

• probabilistic classifier: returns a probability for every label \equiv posterior distribution $p(Y=k|X)$
"soft response"

more general than hard classification: can recover decision function by return the most probable label
Vector of probabilities \mathbb{R}^C

$$p(Y=k|X) \Rightarrow f(x) = \arg \max_k p(Y=k|X)$$

quality measured by "calibration": if $p(Y=k|X)=v$, then label k should be correct $v\%$ of the cases
check with $k=Y^*$

if actual accuracy higher $v\%$: "under confident"
lower $v\%$: "over confident" \leftarrow dangerous, often in neural networks

important: calculate confusion matrix or calibration from test set

Threshold classifier: single feature $X_i \in \mathbb{R}$ $\hat{y} = \text{sign}(X - T) = \begin{cases} 1 & \text{if } X > T \\ -1 & \text{if } X < T \end{cases}$

ex: $y = \begin{cases} 1 & \text{"obese"} \\ -1 & \text{"not obese"} \end{cases}$ $X = \text{weight}$ $T: 90 \text{ kg}$

usually, a single feature is not enough to predict response

three classical solutions for multiple features:

- (1) design a formula to combine feature into single score ex. body-mass index $X_i = \{\text{weight}, \text{height}\}$
but: hard & expensive for most problems $\text{BMI} = \text{weight} / \text{height}^2$ $y = \text{sign}(\text{BMI} - 25)$
- (2) compute score as a linear combination and learn coefficients $S_i = \sum_{j=1}^D \beta_j X_{ij} \Rightarrow \hat{y}_i = \text{sign}(S_i - T)$
"linear classification"
- (3) classify test instance X_{test} according to most similar training instance $\hat{i} = \arg \min_i \text{dist}(X_i, X_{\text{test}})$
"nearest neighbor classification" $\hat{y}_{\text{test}} = y_{\hat{i}}^*$

problems of NN methods: - expensive when N is large: store entire TS and scan for every X_{test}
(faster search exists, but scales badly with D)

- very hard to define $\text{dist}(X_{\text{test}}, X_i)$ to reflect true semantic similarity
 \Rightarrow machine learning task "metric learning" figure out dist from training data

Bayes rule of conditional probabilities:

joint distribution of features & labels: $p(X, Y)$

can be decomposed by chain rule of probability in two ways

posterior \swarrow likelihood \searrow prior

$$p(X, Y) = p(X) p(Y|X)$$

first measure features, then resp

$$\text{classification: } p(Y=k|X) = \frac{p(X|Y=k) p(Y=k)}{p(X)}$$

$$= p(Y) p(X|Y)$$

first determine response then get compatible features

Bayes' rule

$$p(X) = \sum_{k=1}^C p(Y=k) p(X|Y=k)$$

marginal evidence

Bayes Rule: $p(Y=k|X) = \frac{p(X|Y=k) p(Y=k)}{p(X)}$

$p(Y=k|X)$ is labeled **posterior**
 $p(X|Y=k)$ is labeled **likelihood**
 $p(Y=k)$ is labeled **prior**
 $p(X)$ is labeled **evidence / marginal**

$$p(X) = \sum_{k=1}^c p(X|Y=k) \cdot p(Y=k)$$

why is this important?

- general: fundamental equation for probabilistic machine learning because it allows clean uncertainty handling
- defines two fundamental model types
- puts model accuracy into perspective: what is good or bad perf.?

fundamental model types

- discriminative models: learn $p(Y=k|X)$ (LHS of Bayes)
 - answer question "what class" directly
 - ⊕ relatively easy - take direct route, no detour
 - ⊖ often hard to interpret, how model makes decisions
"black box" behavior of neural networks

- generative model: learns $p(Y=k)$ and $p(X|Y=k)$ (RHS of Bayes)
 - first learn, how "the world" works: understand mechanism, then use this to answer question } how observations ("phenotypes") arise from hidden properties ("genes")
 - ⊖ more difficult: need powerful models & a lot of data
 - ⊕ more interpretable

• history:

- traditional science seeks generative models: can create synthetic data that are indistinguishable from real data
- ≈ 1990 : ML researchers realized that their models were too weak \Rightarrow field focused on discriminative models
- ≈ 2012 : neural networks solved many hard discriminative tasks \Rightarrow field is again interested in generative models (ChatGPT, Midjourney) subfield "explainable / interpretable ML"

How does Bayes help to evaluate accuracy?

- how good can a learned classifier be?

Def: Bayes classifier uses Bayes rule (LHS or RHS) with true probs. p^*

Thm: No learned classifier using \hat{p} can be better than Bayes with p^*

ex: $p^*(Y=1|X)=0.6$ then error rate less 1-60% = 40% impossible

how bad can a classifier be? $1 - p^*(Y=Y^*|X)$

case 1: all classes are equally probable $p(Y=k) = \frac{1}{C}$ for all $k=1, \dots, C$
two classes $p(Y=k) = 0.5$

worst classifier: pure guessing \Rightarrow correct 50% of time

case 2: unbalanced classes, e.g. $p(Y=1)=0.01$ $p(Y=-1)=0.99$
 \uparrow prevalence of disease \uparrow does not have this disease
worst classifier: always return majority label \Rightarrow 99% correct

ex: breast cancer screening test, e.g. mammography

$$p(Y=1)=0.01 \quad p(Y=-1)=0.99$$

$$p(X=\text{test positive} | Y=1) = 0.99$$

$$p(X=\text{test positive} | Y=-1) = 0.01$$

"false alarm"

Q: If a test is positive, should you panic?

$$\begin{aligned} p(Y=1 | X=\text{test positive}) &= \frac{p(X=1 | Y=1) p(Y=1)}{p(X=1 | Y=1) p(Y=1) + p(X=1 | Y=-1) p(Y=-1)} \\ &= \frac{0.99 \cdot 0.01}{0.99 \cdot 0.01 + 0.01 \cdot 0.99} = 0.5 \end{aligned}$$