

Classification:

- recap: - regression : Y continuous
 - classification : Y discrete
- without loss of generality : $Y \in \{1, 2, \dots, C\}$ ↑ number of "classes" "labels"
- important special case : two-class classification, $C=2$

$Y \in \{0, 1\}$ or $Y = \{-1, +1\}$
↑ ↑
negative positive outcome
is not infected is infected

- objective: train a function $\hat{Y} = \hat{f}(X)$ such that $\hat{Y} = Y^*$ often ↑ true resp.
 measure the quality: confusion matrix

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

TP rate = $\frac{\# TP}{\# N}$

TN rate = $\frac{\# TN}{N}$

FP rate = $\frac{\# FP}{N}$ FN rate = $\frac{\# FN}{N}$

total error rate = $\frac{\# FP + \# FN}{N}$

should be as small as possible

$N = TP + TN + FP + FN$

- two types of classifiers :
 - hard response : return a single label (hopefully right)
 - soft response : returns the posterior probability $\hat{p}(Y|X)$
 $\hat{p}(Y=k|X)$ is certainty of the classifier that label k is correct

$\hat{p}(Y=k|X) = 0.9 \Rightarrow k$ is probably the true label
 $= 0.1$ is unlikely to be the right label

objective: return well-calibrated probabilities

if $\hat{p}(Y=k|X) = 0.9 \Rightarrow$ the label k should be correct 90% of the time

- under confidence : classifier is right more often
- over confidence : classifier is right less often

Chain rule of probability theory (do not confuse with chain rule in calculus)
considers the joint probability $p(X, Y)$: the current instance has simultaneously the features X and response Y

$$\underset{\substack{\text{observe } X, Y \\ \text{simultaneously}}}{p(X, Y)} = \underset{\substack{\text{observe} \\ X \text{ first}}}{p(X)} \underset{\substack{\text{observe } Y, \\ \text{already knowing } X}}{p(Y|X)} = \underset{\substack{\text{observe } Y \\ \text{first}}}{p(Y)} \underset{\substack{\text{observe } X \\ \text{knowing } Y \text{ already}}}{p(X|Y)}$$
 are all equivalent

$$\begin{aligned} p(X, Y, Z) &= p(X) p(Y|X) p(Z|X, Y) = p(X) p(Z|X) p(Y|X, Z) \\ &= p(Y) p(X|Y) p(Z|X, Y) = p(Y) p(Z|Y) p(X|Y, Z) \\ &= p(Z) p(X|Z) p(Y|X, Z) = p(Z) p(Y|Z) p(X|Y, Z) \end{aligned}$$

in general:
 $M!$ decays
for M variables

Bayes rule: rewrite of the chain rule

$$p(Y|X) = \frac{p(X|Y) \cdot p(Y)}{p(X)}$$

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

$p(Y|X)$: posterior probab. } of Y
 $p(Y)$: prior probability

$p(X|Y)$: likelihood of features for class Y

$p(X)$: evidence (frequency of seeing features X)

normalization: $\sum_k p(Y=k|X) = 1$ for all X

$\Rightarrow p(X)$ normalizes the right-hand side

$$p(X) = \sum_k p(X|Y=k) p(Y=k) \Rightarrow 1 = \sum_k \frac{p(X|Y=k) p(Y=k)}{p(X)}$$

two types of soft classifiers:

- discriminative classifier: learns posterior $\hat{p}(Y|X)$ (LHS of Bayes)
- generative classifier: learns the prior $\hat{p}(Y)$ and the likelihood $\hat{p}(X|Y)$
 (term "generative": if you know the RHS of Bayes, you can create ("generate") synthetic data that are indistinguishable from real data)
- in practice:
 - discriminative cl. are usually more accurate (easier problem)
 - generative cl. are more informative (to create something, you need to understand it, harder)

• worst case performance: features are uninformative about Y

$$p(X|Y=k) = p(X|Y=k') \text{ for all pair } k, k' \in \{1, \dots, C\}$$

$$p(X) = \sum_k p(X|Y=k) p(Y=k) = p(X|Y) \underbrace{\sum_k p(Y=k)}_{=1 \text{ normalization}} = p(X|Y)$$

- which of the two classifiers is preferable to minimize the error probability

error prob: classifier 1: $\hat{f}(x) = 0$ $p^*(Y^* = 1 | X)$

classifier 2: $\hat{f}(x) = 1$ $p^*(Y^* = 0 | X)$

error is minimized when we use cl. 1 if $p^*(Y^* = 1 | X) < p^*(Y^* = 0 | X)$
and cl. 2 otherwise

$$\Leftrightarrow \text{argmax} \quad \hat{f}(x) = \arg \max_{Y_i} p^*(Y^* = Y_i | X)$$

because then the probability of the opposite outcome is minimal

- Bayes classifier: deciding for the most probable outcome for each x
also minimizes the global expected error:

$\mathbb{E}_{x, y^*} [p(\hat{f}(x) \neq y^*)] \rightarrow$ minimized when $p(\hat{f}(x) \neq y^*)$ is

minimized for every x
if classifier has $p^*(Y | X) = p^*(Y | X)$ and returns $\hat{y} = \arg \max_{Y_i} p^*(Y = Y_i | X)$
it achieves the theoretically best possible performance

[catch: in practice, we do not know $p^*(Y | X)$, so we will at best
achieve $p^1 \approx p^*$ \Rightarrow Bayes classifier is a theoretical limit]

- for C labels, the guessing rate to beat is $\frac{1}{C}$

this becomes harder and harder as C increases