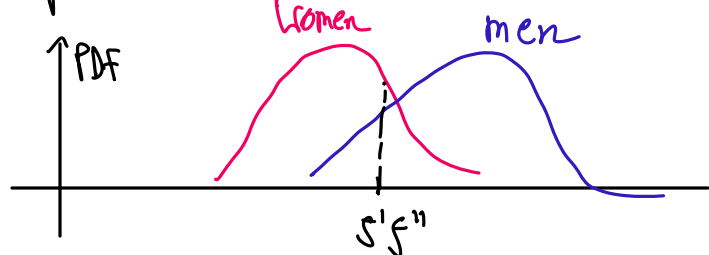# Decision Theory (risk minimization)

adult height $\longrightarrow$ Man/Woman

1st issue: multiple pts w/ different classes could lie at same pt

Want: probabilistic classifier



Suppose 10% of population has cancer, 90% doesn't

Prob distributions for calorie intake, $P(x|y)$: $X = \#$ Calories $Y = $ Cancer

| Calories $(X)$ | $\leq 1,200$ | $1200 - 1600$ | $> 1600$ |
|---|---|---|---|
| Cancer $(y=1)$ | 20% | 50% | 30% |
| $\overline{Cancer}$ $(y=-1)$ | 1% | 10% | 89% |

Recall: $P(x) = P(x|y)P(y) + P(x|\overline{y})P(\overline{y})$

$$P(1200 \leq x \leq 1600) = 0.5 \cdot 0.1 + 0.9 \cdot 0.10 = 0.14$$

guy eats $X = 1400$ Calories. Prob(guy has cancer)

bayes rule = $\underbrace{P}_{\text{(posterior probability)}} \dfrac{\text{Prob}(12-1600 | \text{cancer}) \overbrace{P(Cancer)}^{\text{prior probability}}}{P(1200 \leq x \leq 1600)} = \dfrac{0.2 \cdot 0.1}{0.14} = \dfrac{0.05}{0.14}$

*assumes equal loss for misclassifications
↳ need to punish false negatives more

$\longrightarrow P(\text{cancer} | 1200 \leq x \leq 1600) = \dfrac{5}{14} \approx 36\%$

loss function $L(z,y)$ specifies badness if classifier predicts $z$ but classifies $y$

e.g. $L(z,y) = \begin{cases} 1 & z=1, y=-1 \quad \text{false positive is bad} \\ 5 & z=-1, y=1 \quad \text{false negative really bad} \\ 0 & z=y \quad \text{correct classification good} \end{cases}$

↗ asymmetrical

36% probability of loss 5, worse than 64% chance of loss 1

Definitions:

Symmetric loss fn: penalties are equal

0-1 loss fun: 0 if True, 1 if false/incorrect

let $r: \mathbb{R}^d \to \pm 1$ (decision rule/classifier)

$\quad \hookrightarrow$ maps feature vector to prediction

Risk for a decision rule $R$ = expected loss over all possible values of $x$ & $y$

$$\text{Risk} = \text{expected loss}$$

$\to R(r) = E[L(r(x), y)]$

functional

$$= \sum_x \left[ (L(r(x), 1) \cdot P(y=1|x) + L(r(x), -1) P(y=-1|x=x) \right] P(x=x)$$

outcome — prob of outcome

Posterior

set of all possible values of RV $x$

Points not in class

$$(\text{Baye's rule}) = P(y=1) \sum_x L(r(x), 1) P(x=x|y=1) + P(y=-1) \sum_x L(r(x), -1) P(x=x|y=-1)$$

prior

Bayes decision rule: function $r^*$ that minimizes functional $R(r)$
(Bayes Classifier) (takes in a function)

Assuming $L(z, y) = 0$ for $z = y$:

discrete case $\quad r^*(x) = \begin{cases} 1 & \text{if } \boxed{L(-1,1) P(y=1|x=x) > L(1,-1) P(y=-1|x=x)} \\ -1 & \text{o/w} \end{cases}$

(which of 2 terms is bigger)

When $L$ is symmetric, $L(-1,1) = L(1,-1) \implies$ only compare $P(y=1|x=x) > P(y=-1|x=x)$

$\quad \hookrightarrow$ In other words: Pick class w/ biggest posterior probability

If not symmetric: weigh posteriors by respective loss functions

for cancer example: $r^*(x) = 1$ for $x \leq 1600$, $r^*(x) = -1$ for $x > 1600$

he chose these

Bayes/optimal risk: risk of best possible classifier (Bayes classifier)

Cancer example: $R(r^*) = 0.1(5 \cdot 0.3) + 0.9(1 \times 0.01 + 1 \cdot 0.1) = \underline{0.249}$

*using prior prob. formula

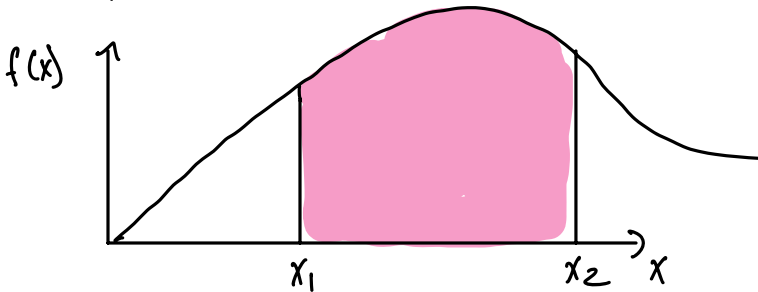no decision rule gives lower risk

Deriving/using $r^*$ is called risk minimization

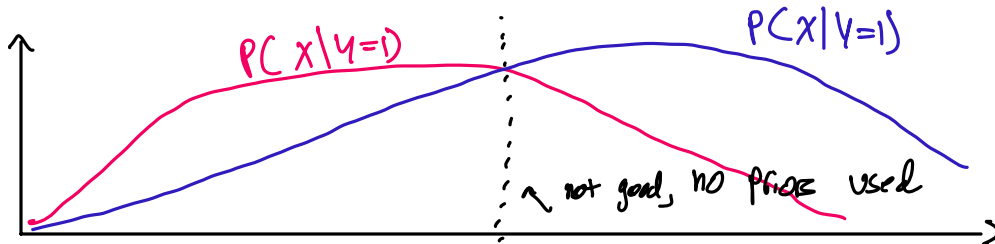# CONTINUOUS DISTRIBUTIONS

$X$ is continuous w/ a PDF

review:



prob that $X \in [x_1, x_2] = \int_{x_1}^{x_2} f(x) dx$

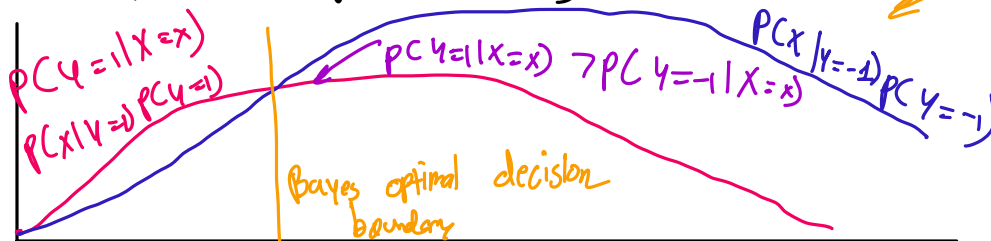$X \in [-\infty, \infty] = 1$ (def of pdf) $= \int_{-\infty}^{\infty} f(x) dx = \underline{1}$

Expected value of $g(x) = \int_{-\infty}^{\infty} g(x) f(x) dx$

mean $\mu = E[x] = \int_{-\infty}^{\infty} x f(x) dx$

Variance: $\sigma^2 = E[(x-\mu)^2] = E[x^2] - E[x]^2$



$P(X|Y=-1)$

$P(X|Y=1)$

↖ not good, no priors used

Suppose $P(Y=1) = 1/3$, $P(Y=-1) = 2/3$, 0-1 loss

$P(Y=1|X=x)$

$P(X|Y=1) P(Y=1)$

$P(Y=1|X=x) > P(Y=-1|X=x)$

$P(X|Y=-1) P(Y=-1)$

Posterior probabilities, don't need to divide by $P(X)$

Bayes optimal decision boundary

Define risk as before, Integral instead of summation

$$R(r) = E\left[L(r(x), y)\right]$$

$$= P(Y=1) \int L(r(x), 1) f(X=x | Y=1) dx$$

$$+ P(Y=-1) \int L(r(x), -1) f(X=x | Y=-1) dx$$

Bayes risk = area under minimum of 2 functions

Assuming $L(z, y) = 0$ for $z = y$:

$$R(r^*) = \int \min_{y=\pm 1} L(-y, y) f(X=x | Y=y) P(Y=y) dx$$

0-1 loss $\Longrightarrow$ Risk $= P(r(x)$ is wrong$)$

posterior probability

$\hookrightarrow$ Bayes optimal decision boundary $= \{ x : \underbrace{P(Y=1 | X=x)}_{\text{decision function}} = \underbrace{0.5}_{\text{isovalue}} \}$

## 3 WAYS TO BUILD Classifiers

① generative models ( e.g. LDA)
  - Assume sample pts come from probability distributions, different for each class
  - guess form of distributions
  - For each class $C_i$, fit distribution parameters to the class, giving $f(x | Y=C_i)$
  - estimate prior probability for each class $P(Y=c)$
  - get posterior probability w/ Bayes thm
  - If 0-1, pick $c$ that maximizes $P(Y=c | X=x)$ equivalent to maximizing $f(X=x | Y=c) P(Y=c)$

② Discriminative Models ( e.g. logistic regression)
  - model posterior probability w/o class conditional probabilities

③ Find decision boundary ( e.g. SVM)

  - model $r(x)$ directly, (no posterior)

Advantages of 1 & 2 :- $P(Y|x)$ tells probability you're wrong

advantage of ①: easy to diagnose outliers : $f(x)$ very small

disadvantage of ①: hard to estimate distributions, e.g. real distributions