# Probability & Inference

- Result of tossing a coin is $\in \{H, T\}$

- Random variable $X \in \{1, 0\}$

  — Recall the Bernoulli PMF

  $$P(X=1) = \theta^x (1-\theta)^{(1-x)}$$

- Sample: $X = \{X^t\}_{t=1}^N$

- Estimation (MLE): $\hat{\theta} = \dfrac{\# H}{\# Tosses} = \dfrac{1}{N} \sum_{t=1}^{N} X^t$

- Predict the next toss:

  $$\begin{cases} H & \text{if } \hat{\theta} > \frac{1}{2}, \\ T & \text{otherwise} \end{cases}$$

Conjugate Bayes
----
$\theta \sim Beta(a, b)$

$\theta | x \sim Beta(a^*, b^*)$

$a^* = a + N_1$

$b^* = b + N_0$

$$\left\{ \text{where } \begin{array}{l} N_1 = \# \text{ heads} \\ N_0 = \# \text{ tails} \\ N_0 + N_1 = N \# \text{ of trials} \end{array} \right.$$

# Classification

- Credit scoring: Inputs are income and savings

  — Output: low-risk vs high-risk

- Input: $X = [x_1, x_2]^T$

  Output: $C \in \{0, 1\}$

- Prediction

  choose/select $\begin{cases} C=1 & \text{if } P(c=1 \mid x_1, x_2) > 0.5 \\ C=0 & \text{otherwise} \end{cases}$

  or

  choose/select $\begin{cases} C=1 & \text{if } P(c=1 \mid x_1, x_2) > P(c=0 \mid x_1, x_2) \\ C=0 & \text{otherwise} \end{cases}$
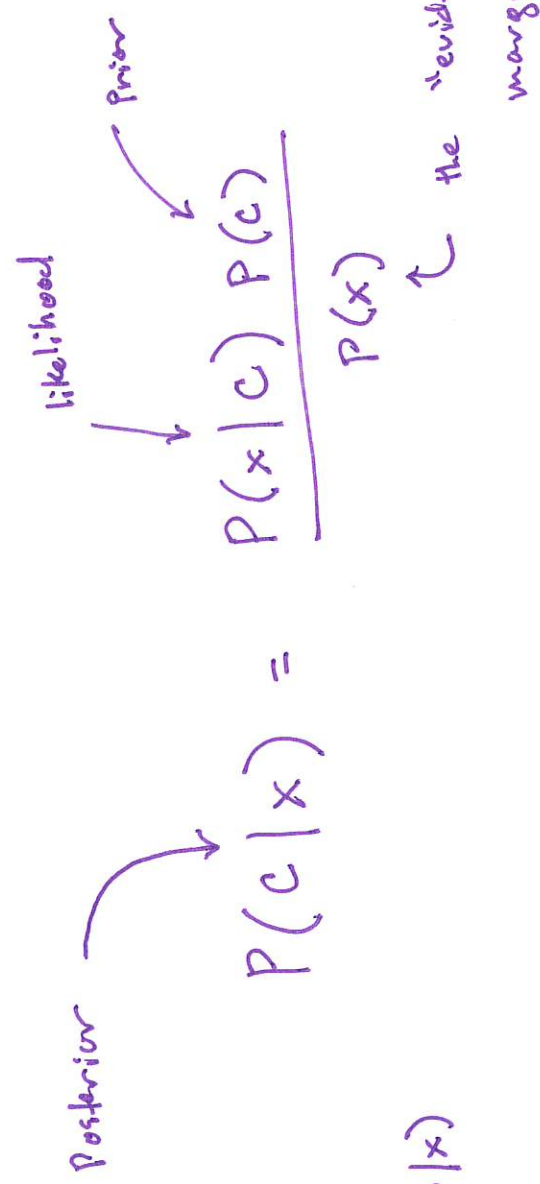
# Bayes' Rule

(3)

NB:
$$\begin{cases} P(c=0|x) + P(c=1|x) = 1 \end{cases}$$

## Decision

choose $C = 1$ if

$P(c=1|x) > P(c=0|x)$

Posterior $\longrightarrow$

likelihood

Prior

$$P(c|x) = \frac{P(x|c) \, P(c)}{P(x)}$$

the "evidence" or

marginal of the data

Prior prob. $P(C=1)$ : prob. of high-risk customer

$- \quad P(c=0) + P(C=1) = 1$

Class Likelihood : $P(x|c)$    prob. of event belonging to class $C$ given

the observable $x$   $\left(\begin{array}{l} Q: \text{how } \text{likely} \\ \text{are the data if the true} \\ \text{state of the world is } C \end{array}\right)$

$- \quad P(x_1, x_2 | c) =$ prob. that high-risk customer

will have $[x_1, x_2]^T$

Evidence : $P(x)$ : marginal prob. of $[x_1, x_2]^T$ $\begin{cases} P(x) = P(x|c=1) P(c=1) + P(x|c=0)P(c=0) \end{cases}$

④

Bayes' Rule for $k > 2$ Classes    ( mutually exclusive & exhaustive )

Posterior prob

$\downarrow$

$$P(c_i | x) = \frac{P(x | c_i) \, P(c_i)}{P(x)}$$

$$= \frac{P(x | c_i) \, P(c_i)}{\sum_{k=1}^{k} P(x | c_k) \, P(c_k)}$$

Prior Probabilities :  $P(c_i) \quad \forall \; i = 1, \ldots, k$

$$\sum_{i=1}^{k} P(c_i) = 1$$

Class Likelihood :  $P(x | c_i)$

Decision:

Choose as the output the class $c_i$

that has the maximum posterior prob.

$$P(c_i | x)$$

# Losses and Risks

- Credit scoring decisions should be made so as to maximize gains / limit losses

- Classes $c_1, \ldots, c_k$

- Let $d_i$ = decision to assign $c_i$ to the input $\forall$ $1 \leq i \leq k$

- Let $\lambda_{ik}$ = the loss from assigning $c_i$ to input that belongs to $c_k$ (i.e. misclassification)

- Expected risk for $d_i$:

$$R(d_i \mid x) = \sum_{k=1}^{k} \lambda_{ik} \, P(c_k \mid x)$$

- Decision: choose $d_i$ with minimal expected risk

$$R(d_i \mid x) = \min_k R(d_k \mid x)$$

# 0/1 Loss Case

- Correct decisions have $\emptyset$ loss

  Incorrect decisions have 1 loss

$$\lambda_{ik} = \begin{cases} \emptyset & \text{if } i = k \\ 1 & \text{if } i \neq k \end{cases}$$

NB:
- all equal
- losses are symmetric

- The risk of taking action $\alpha_i$ is:

$$R(\alpha_i \,|\, x) = \sum_{k=1}^{K} \lambda_{ik} P(c_k \,|\, x) = \sum_{\substack{k=1 \\ k \neq i}}^{K} P(c_k \,|\, x) = 1 - P(c_i \,|\, x)$$

$\uparrow$ 
$\emptyset$ or 1

via compliment rule

- Decision: to minimize risk, assign $c_i$ to the most probable class

# Losses & Risks: Reject

- reject option might require human analysts
- or another ML algorithm

- misclassification: very high cost

- Consider an extra option $(K+1)$-st class "reject"

- Loss function:

$$\lambda_{ik} = \begin{cases} 0 & \text{if } i = k \\ \lambda & \text{if } i = k+1, \quad 0 < \lambda < 1 \\ 1 & \text{if } i \neq k \text{ and } i = k+1 \end{cases}$$

- Risk of reject:

$$R(d_{K+1} | x) = \sum_{k=1}^{K} \lambda\, P(C_k | x) = \lambda$$

- Risk of misclassification:

$$R(d_i | x) = \sum_{k \neq i}^{l} P(C_k | x) = 1 - P(C_i | x)$$

via compliment rule (i.e. one minus "doing it correctly")

## Decision:

- Choose $C_i$; $\min_{1 \leq i \leq K} R(d_i | x)$

output:
- $C_i$ if $P(C_i | x) > P(C_k | x)$
- $C_i$ if $P(C_i | x) > 1 - \lambda$
- Reject otherwise

## Cases:

$\lambda = 0$: "always reject" (as good as correct classification)

$\lambda = 1$: "never reject" (as bad as incorrect classification)

## NB:

$$\lambda \underbrace{\sum_{k=1}^{K} P(C_k | x)}_{1}$$

# Discriminant Functions

- Choose $g_i(x)$, $i = 1, \ldots, k$  s.t. output $C_i$  if  $g_i(x) = \max_k g_k(x)$
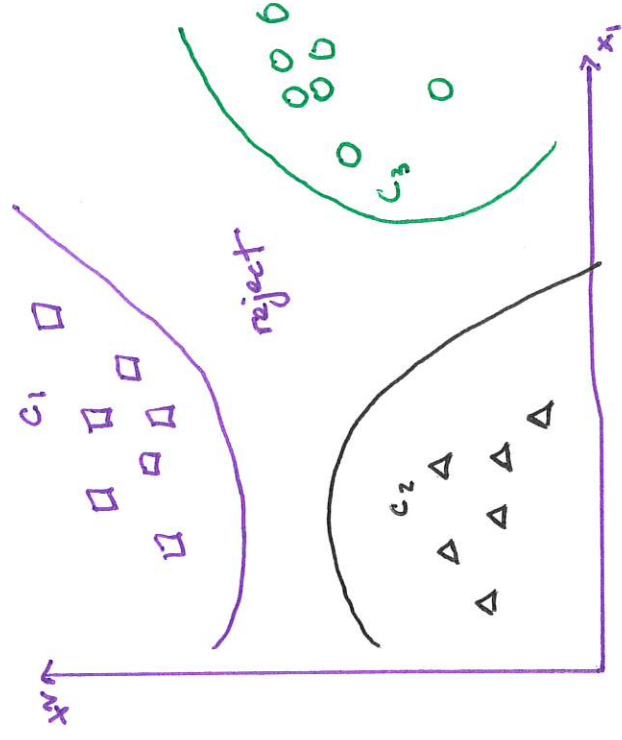
- **Bayes' Classifier:**

$$g_i(x) = -R(\alpha_i | x)$$

  — 0/1 loss: $g_i(x) = P(C_i | x) = \frac{P(x | c_i) \, P(C_i)}{P(x)}$  via Bayes' Rule

  Equivalent way: $g_i(x) = P(x | c_i) \, P(c_i)$  (simpler / easier)

- $k$ decision regions

$$R_i = \left\{ x \;\middle|\; g_i(x) = g_k(x) \right\}$$

# $k = 2$ Classes

- Dichotomizer $(k=2)$ vs Polychotomizer $(k>2)$

- $g(x) = g_1(x) - g_2(x)$

choose $\begin{cases} c_1 & \text{if } g(x) > 0 \\ c_2 & \text{otherwise} \end{cases}$

- Log odds:

$$\log \frac{P(c_1 \mid x)}{P(c_2 \mid x)} = \underbrace{\log P(c_1 \mid x) - \log P(c_2 \mid x)}_{\text{log difference in probabilities}}$$

# Bayesian Thinking Revisited

Introduction to Bayes Factors

NB: notes based on chp. 3 Losses and Decision Making

An Introduction to Bayesian Thinking

- Return to HIV testing

  - $H_1$: Patient does not have HIV

  - $H_2$: Patient does have HIV

- Actions / Decisions

  - $d_1$: Choose $H_1$

  - $d_2$: choose $H_2$

  - $L(d)$ the loss associated w/ decision $d_i$     $\forall i \in \{1, 2\}$

    - $d = d_1$
      - Right: decide $H_1$ / and they don't   $\Rightarrow L(d_1) = 0$
      - Wrong: decide $H_1$ / and they do   $\Rightarrow L(d_1) = w_1$

    - $d = d_2$
      - Right: decide $H_2$ / they do $\Rightarrow L(d_2) = 0$
      - Wrong: decide $H_2$ / they don't $\Rightarrow L(d_2) = w_2$

## Losses

$H_1:$ Patient has HIV ~~does not~~ (does not)

$H_2:$ Patient has HIV

$$L(d_1) = \begin{cases} 0 & \text{if } d_1 \text{ correct} \\ w_1 = 1000 & \text{else} \end{cases}$$

$$L(d_2) = \begin{cases} 0 & \text{if } d_2 \text{ correct} \\ w_2 = 10 & \text{else} \end{cases}$$

## Posteriors

- $(+)$ stands for a positive result from the ELISA

- $P(H_1 \mid +) \approx 0.88$

    posterior prob. of NOT having HIV given
    a $(+)$ ELISA result

- $P(H_2 \mid +) \approx 0.12$

    Posterior prob. of having HIV given the
    $(+)$ ELISA result

    (from the complement)

Expected Losses

- $E(L(d_1)) = 0.88(0) + 0.12(1000) = 120$

  $E[L(d_2)] = 0.88(10) + 0.12(0) = 8.8$

- Since $E[L(d_2)] < E[L(d_1)]$ $\Rightarrow$ decide the patient has HIV

NB:
- Decision highly influenced by losses assigned to $d_1$ and $d_2$

- If losses symmetric, say $\omega_1 = \omega_2 = 10$

  $- E[L(d_1)] = 0.88(0) + 0.12(10) = 1.2$

  $-$ While $E[L(d_2)]$ would not change

  $-$ We would decide that patient does NOT have HIV!

# Bayes Factors

– Continue with HIV testing example

– Prior odds = ratio of the prior probabilities of hypotheses

$$O[H_1 : H_2] = \frac{P(H_1)}{P(H_2)}$$

– Posterior odds = ratio of the two posterior probabilities
of hypotheses

$$PO[H_1 : H_2] = \frac{P(H_1 | Data)}{P(H_2 | Data)}$$

Using Bayes' Rule, we find that

$$PO[H_1 : H_2] = \frac{P(H_1 | Data)}{P(H_2 | Data)}$$

$$= \frac{P(data | H_1) P(H_1) \Big/ P(data)}{P(data | H_2) P(H_2) \Big/ P(data)}$$

$$= \frac{P(data | H_1) \; P(H_1)}{P(data | H_2) \; P(H_2)}$$

$$= \underbrace{\frac{P(data | H_1)}{P(data | H_2)}}_{\text{Bayes Factor}} * \underbrace{\frac{P(H_1)}{P(H_2)}}_{\text{Prior odds}}$$

So, we have that

$$PO[H_1 : H_2] = BF[H_1 : H_2] * O[H_1 : H_2]$$

– BF quantifies the evidence of data arising

from $H_1$ versus $H_2$

– In the discrete case:

$$BF[H_1 : H_2] = \frac{P(data \mid H_1)}{P(data \mid H_2)}$$

(likelihood ratio)

– In the continuous case

$$BF [H_1 : H_2] = \frac{\int P(data \mid \theta, H_1) \, d\theta}{\int P(data \mid \theta, H_2) \, d\theta}$$

NB: $\theta$ is the index of all possible models / parameters

— Continue w/ HIV case

  $H_1$: Patient does not have HIV

  $H_2$: Patient does have HIV

— Priors:

  $P(H_1) = .99852$

  $P(H_2) = .00148$

NB: from prevalence
of HIV in the
population

— Prior odds then is

  $$O[H_1 : H_2] = \frac{P(H_1)}{P(H_2)} = \frac{.99852}{.00148} = 674.6757$$

- Posteriors

$$P(H_1 \mid +) = .878855$$

$$P(H_2 \mid +) = .121144$$

$$PO[H_1 : H_2] = \frac{.121144}{.878855} = 2.254578$$

- Bayes Factor

$$BF[H_1 : H_2] = \frac{O[H_1 : H_2]}{PO[H_1 : H_2]}$$

$$= \frac{.125454}{674.9457} = 0.0108$$

$$= \frac{P(+ \mid H_1)}{P(+ \mid H_2)} = \frac{.01}{.93} \approx .0108$$

- So now that we have calculated $BF[H_1:H_2]$, how should understand it's meaning?

- Jeffrey (1961)

| $BF[H_1:H_2]$ | Evidence against $H_2$ |
|---|---|
| 1 to 3 | Not worth a bare mention |
| 3 to 20 | Positive |
| 20 to 150 | Strong |
| > 150 | Very strong |

| $2 * \log \left( BF[H_1 : H_2] \right)$ | Evidence against $H_1$ |
|---|---|
| 0 to 2 | Not worth a bare mention |
| 2 to 6 | Positive |
| 6 to 10 | Strong |
| > 10 | Very strong |

- NB: notice that for the HW case it doesn't even appear on the scale.

- Let's reader so that it's $BF[M_2 : H_1] = \dfrac{1}{BF[H_1 : H_2]} = \dfrac{1}{.01085} = 92.54259$

$$\approx 93$$

- Hence the evidence against $H_1$ is "very strong"