

# M2. Bayesian Decision Theory (incl. Bayes classifiers)

Manikandan Narayanan

Feb 4-5,10-11,15 2021

PRML Jan-May 2021 (MKN section)

# Acknowledgment of Sources

- Slides based on content from related
  - Courses:
    - IITM – Profs. Arun/Harish/Chandra’s PRML offerings (slides, quizzes, notes, etc.), Prof. Ravi’s “Intro to ML” slides – cited respectively as [AR], [HR], [CC], [BR] in the bottom right of a slide.
    - India – NPTEL PR course by IISc Prof. PS. Sastry (slides, etc.) – cited as [PSS] in the bottom right of a slide.
  - Books:
    - PRML by Bishop. (content, figures, slides, etc.) – cited as [CMB]
    - Pattern Classification by Duda, Hart and Stork. (content, figures, etc.) – [DHS]
    - Mathematics for ML by Deisenroth, Faisal and Ong. (content, figures, etc.) – [DFO]

# Outline of Module M2

- M2. Bayesian Decision Theory (incl. Bayes classifiers)
  - M2.0 Introduction/Background (on Probability Theory)
  - **M2.1 Bayesian Decision Theory**
    - **M2.1.0 Decision Theory for Classification/Regression (common defns./notations)**
    - M2.1.1 Decision Theory for Classification (Bayes classifiers)
    - M2.1.2 Decision Theory for Regression (Squared loss, etc.)

# M2.1.0 Decision Theory (for classification/regression)

$x$  is feature vector (input),  $t$  is target/response (output).

- Inference step
  - Determine either  $p(t|x)$  or  $p(x, t)$ .
- Decision step
  - For any given  $x$ , determine optimal  $t$ .
  - Optimality wrt *risk* or *expected loss*; General loss functions are:
    - Classification ( $t$  discrete): ***misclassification rate***, loss-matrix based function, etc.
    - Regression ( $t$  continuous): ***squared loss***, Minkowski loss, etc.

# Notations

- Feature vector  $\mathbf{x} \in \mathcal{X}$ 
  - Feature vector  $\mathbf{x} = (x_1, x_2, \dots, x_D)$
  - Feature space  $\mathcal{X} = \mathbb{R}^D$ 
    - Think of  $D=1$  in rest of slides, but Bayesian decision theory (Bayes classifier) holds for any  $D$ .
- Target/response  $t \in \mathcal{Y}$ 
  - Discrete: Target space  $\mathcal{Y} = \{C_1, C_2, \dots, C_K\}$ 
    - Often times also referred to as  $\{1, 2, \dots, K\}$ , or for binary ( $K=2$ ) classifiers as  $\{0, 1\}$  or  $\{-1, +1\}$
  - Continuous: Target space  $\mathcal{Y} = \mathbb{R}$
- Classifier or regressor is simply a function from feature to target space
  - i.e., it maps each point in the feature space to a unique point in the target space
  - $h: \mathcal{X} \rightarrow \{C_1, \dots, C_K\}$
  - $y: \mathcal{X} \rightarrow \mathbb{R}$

# Notations (Bayes rule)

- $$P(t|x) = \frac{P(t)P(x|t)}{P(x)} \propto P(t)P(x|t)$$

(posterior = prior x likelihood (class conditional) / evidence)

- $$P(x, t) = P(x) P(t|x) = P(t) P(x|t)$$

(joint = evidence x posterior = prior x liklhd. (class cond.))

- For binary  $t$ , 
$$\begin{aligned} P(x) &= P(t = C_1)P(x|C_1) + p(t = C_2)P(x|C_2) \\ &= P(C_1)P(x|C_1) + P(C_2)P(x|C_2) \end{aligned}$$

# Outline of Module M2

- M2. Bayesian Decision Theory (incl. Bayes classifiers)
  - M2.0 Introduction/Background (on Probability Theory)
  - **M2.1 Bayesian Decision Theory**
    - M2.1.0 Decision Theory for Classification/Regression (common defns./notations)
    - **M2.1.1 Decision Theory for Classification (Bayes classifiers)**
    - M2.1.2 Decision Theory for Regression (Squared loss, etc.)

# M2.1.0 Decision Theory for Classification

- Inference step
  - Determine either  $p(x, t)$  or  $p(t = C_k | x)$ .
- Decision step
  - For any given  $x$ , determine optimal class label  $h(x) = C_j$  for  $t$ .
  - Optimality wrt *risk* or *expected loss* (misclassification rate or general loss function/matrix for binary vs. multi-class classifiers)



# Bayes classifier (two classes)

- $h(\mathbf{x}) = C_1$  if  $P(C_1|\mathbf{x}) > P(C_2|\mathbf{x})$   
     $= C_2$  o. w (otherwise i. e.,  $P(C_2|\mathbf{x}) \geq P(C_1|\mathbf{x})$ )

(Note:       $P(C_1|\mathbf{x}) > P(C_2|\mathbf{x}) \Leftrightarrow$                        $\leftarrow$  for discriminative models  
               $P(C_1, \mathbf{x}) > P(C_2, \mathbf{x}) \Leftrightarrow$                        $\leftarrow$  for generative models  
               $P(C_1)P(\mathbf{x}|C_1) > P(C_2)P(\mathbf{x}|C_2)$   $\leftarrow$  for gen. models' learning)

- Bayes classifier is the ***optimal*** classifier among all classifiers
  - wrt minimizing the probability of error (aka misclassification rate), ...
  - ...assuming complete knowledge of the posterior distribution.

# Optimality – minimum misclassification rate

- Let Decision region  $R_i := \{x \in \mathcal{X} \mid h(x) = C_i\}$

# Optimality – minimum misclassification rate

- Let Decision region  $R_i := \{x \in \mathcal{X} \mid h(x) = C_i\}$

$$P(\text{error}) = P_{x,t}(h(x) \neq t)$$

# Optimality – minimum misclassification rate

- Let Decision region  $R_i := \{x \in \mathcal{X} \mid h(x) = C_i\}$

$$P(\text{error}) = P_{x,t}(h(x) \neq t)$$

$$= \int \sum_{t=C_1, C_2} P(x, t) \mathbb{1}_{\{h(x) \neq t\}} dx$$

# Optimality – minimum misclassification rate

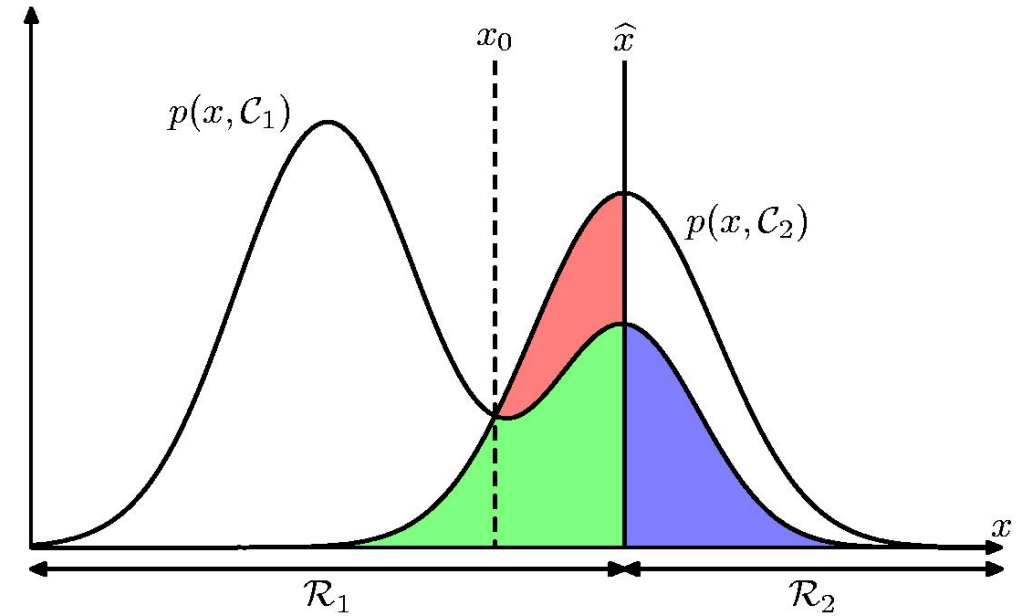
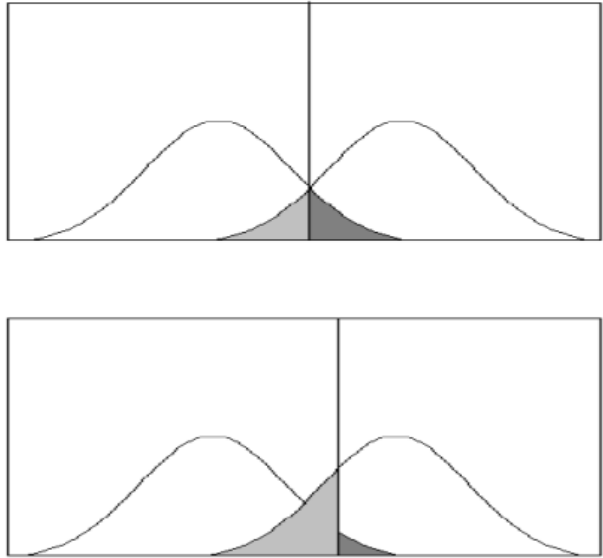
- Let Decision region  $R_i := \{x \in \mathcal{X} \mid h(x) = C_i\}$

$$P(\text{error}) = P(h(x) \neq t)$$

$$= \int \sum_{t=C_1, C_2} P(x, t) \mathbb{1}_{\{h(x) \neq t\}} dx$$

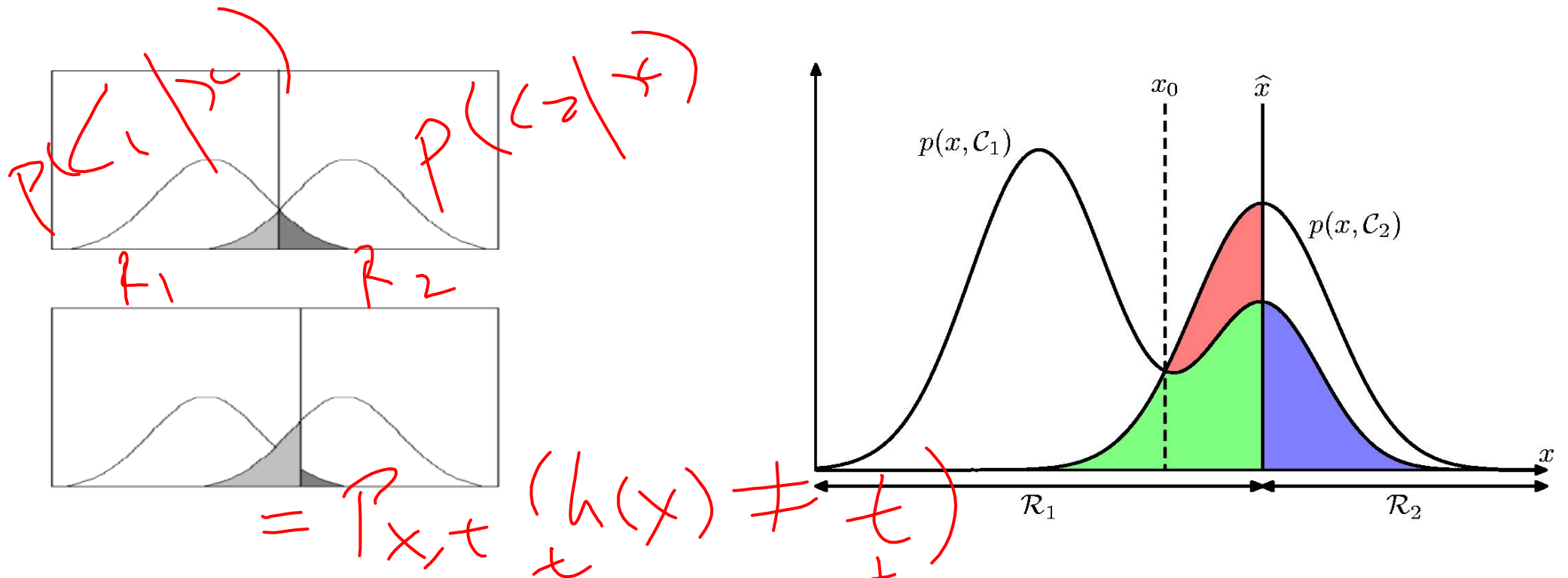
$$= \int \sum_t P(t|x) \mathbb{1}_{\{h(x) \neq t\}} P(x) dx$$

# Optimality - minimum misclassification rate



$$\begin{aligned} p(\text{mistake}) &= p(\mathbf{x} \in \mathcal{R}_1, \mathcal{C}_2) + p(\mathbf{x} \in \mathcal{R}_2, \mathcal{C}_1) \\ &= \int_{\mathcal{R}_1} p(\mathbf{x}, \mathcal{C}_2) d\mathbf{x} + \int_{\mathcal{R}_2} p(\mathbf{x}, \mathcal{C}_1) d\mathbf{x}. \end{aligned}$$

# Optimality - minimum misclassification rate

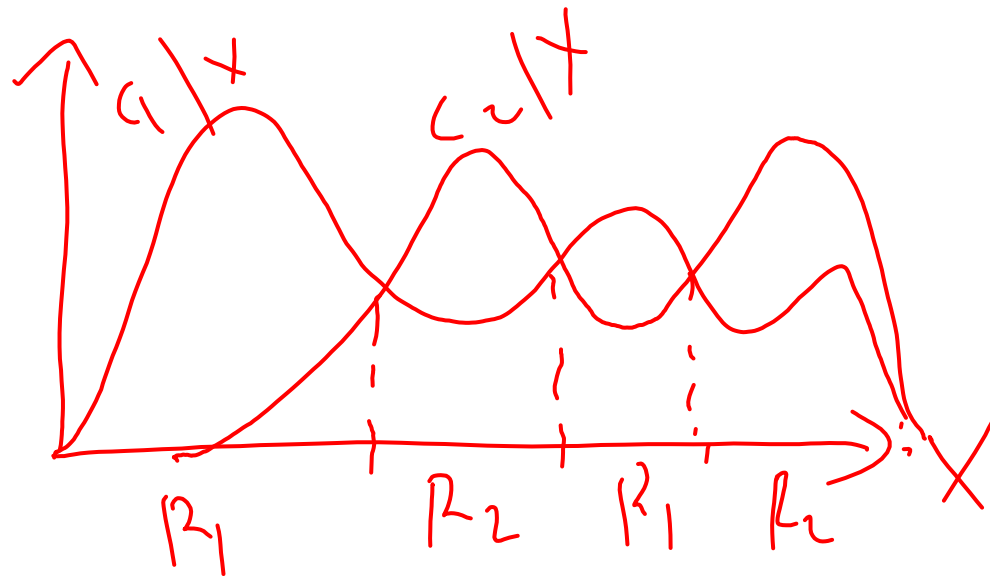


$$\begin{aligned}
 p(\text{mistake}) &= p(\mathbf{x} \in \mathcal{R}_1, \mathcal{C}_2) + p(\mathbf{x} \in \mathcal{R}_2, \mathcal{C}_1) \\
 &= \int_{\mathcal{R}_1} p(\mathbf{x}, \mathcal{C}_2) d\mathbf{x} + \int_{\mathcal{R}_2} p(\mathbf{x}, \mathcal{C}_1) d\mathbf{x}.
 \end{aligned}$$

Can decision regions be discontinuous in the optimal classifier?



Can decision regions be discontinuous in the optimal classifier?



# Bayes classifier (multi-class; $K > 2$ classes)

- $h(x) = C_j$  if  $P(t = C_j | x) \geq P(t = C_{j'}, | x) \quad \forall j' \in \{1, \dots, K\} \setminus \{j\}$   
=  $\operatorname{argmax}_{C_j} P(t = C_j | x)$  (ties broken arbitrarily)

$$h(x) = C_1 \text{ if } p(C_1|x) > p(C_2|x)$$

- Again **optimal** classifier among all classifiers
  - wrt same criteria as for binary classifier i.e., minimum misclassification rate (or) equivalently maximum classification accuracy...
  - ...assuming complete knowledge of the posterior distribution

# Optimality of multi-class classifier (max. accuracy)

optimal  $h$

$$p(x_1) p(h(x_1) | x_1) \\ + p(x_2) p(h(x_2) | x_2)$$

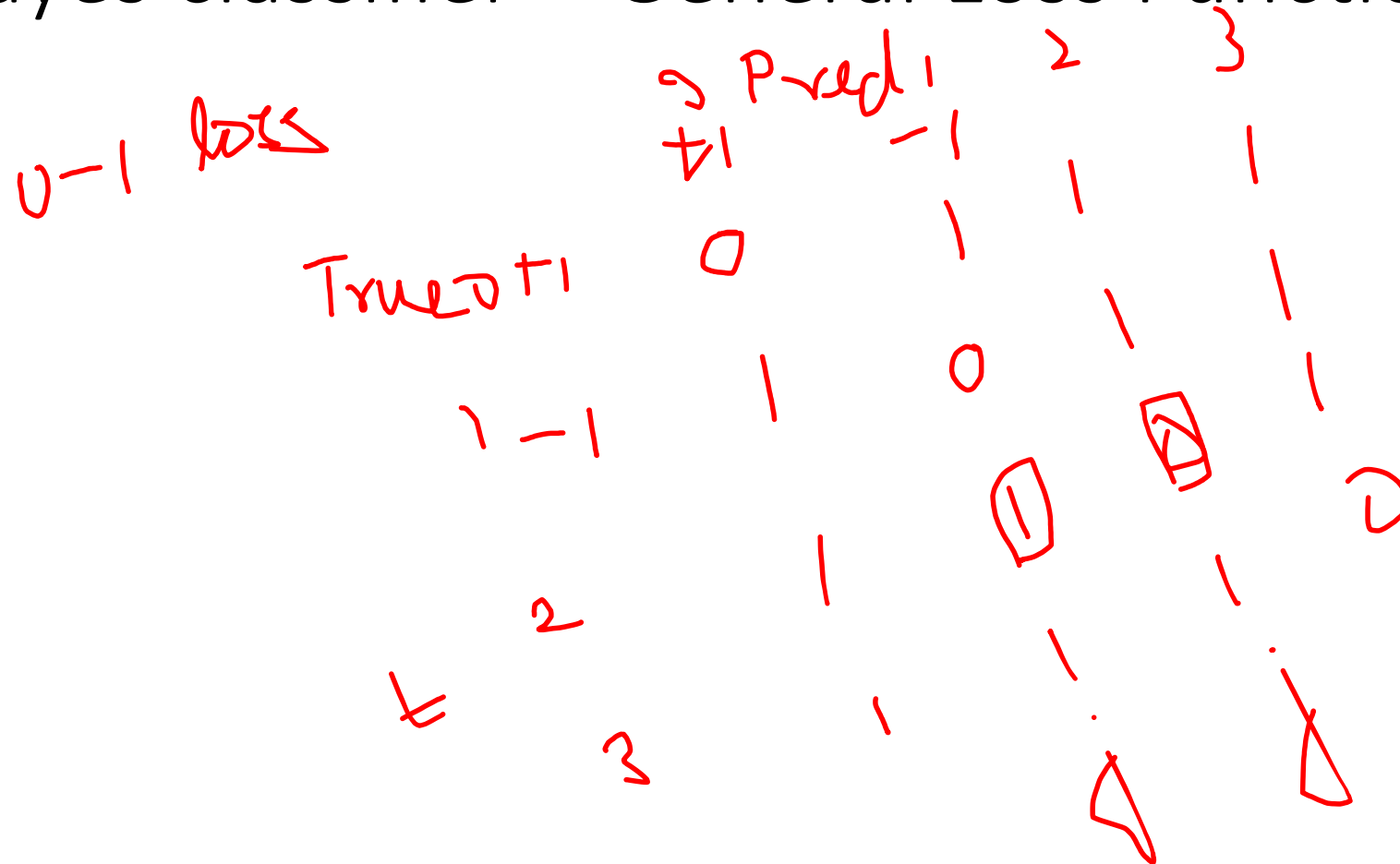
$$p(\text{correct}) = \sum_{k=1}^K p(\mathbf{x} \in \mathcal{R}_k, \mathcal{C}_k) \\ = \sum_{k=1}^K \int_{\mathcal{R}_k} p(\mathbf{x}, \mathcal{C}_k) d\mathbf{x}$$



$$h(x_1) = +1 \\ h(x_2) = +1 \\ h(x_3) = -1 \\ h(x_4) = -1$$

$$p(+1 | x_1) = 0.7 \\ p(+1 | x_2) = 0.4 \\ p(+1 | x_3) = 0.2 \\ p(+1 | x_4) = 0.6$$

# Bayes classifier – General Loss Function (Matrix)



# Bayes classifier – General Loss Function (Matrix)

- Example: classify medical images as 'cancer' or 'normal'

		Decision	
		cancer	normal
Truth	cancer	0	<u>1000</u>
	normal	1	0

# Optimality - Minimum Expected Loss

$$\mathbb{E}[L] = \sum_{t=C_1, \dots, C_K} \sum_{j=C_1, \dots, C_K} \int_{\mathcal{R}_j} L_{tj} p(x, t) dx$$

Regions  $\mathcal{R}_j$  are chosen to minimize (next slides show why)

$$\mathbb{E}[L \mid X = x] = \sum_{t=C_1, \dots, C_K} L_{t, h(x)=j} p(t|x)$$

# Optimality - Minimum Expected Loss

$$\mathbb{E}[L] = \sum_{t=C_1, \dots, C_K} \sum_{j=C_1, \dots, C_K} \int_{\mathcal{R}_j} L_{tj} p(x, t) dx$$

Regions  $\mathcal{R}_j$  are chosen to minimize (next slides show why)

$$\mathbb{E}[L | X = x] = \sum_{t=C_1, \dots, C_K} L_{t, h(x)=j} p(t|x)$$

$L = 0-1$  loss  
 $h(x) = \arg \max_j p(t=j|x)$

$$\equiv h(x) = \arg \min_j \sum_t L_{tj} p(t|x)$$

$\sum_t L_{tj} p(t|x)$   
 $1 \times p(C_2|x)$

# Optimality - Minimum Expected Loss (indicator fn. notation)

$$E[L] = \int \sum_{t \in C_1}^{C_2} p(t|x) \left[ \sum_{j \in C_1}^{C_2} L_{tj} \mathbb{1}_{\{h(x)=j\}} \right] p(x) dx$$

$$= \int \sum_j \left( \underbrace{\sum_t L_{tj} p(t|x)}_{\substack{\text{Choose} \\ h(x)=j \text{ s.t. } \Rightarrow \text{is minimized}}} \right) \mathbb{1}_{\{h(x)=j\}} p(x) dx$$



# Optimality - Minimum Expected Loss (cond. expectation notation)

$$\begin{aligned} E_{x,t}[L] &= E_x[E_{t|x}[L]] \quad (E_x[E_t[L|x]]) \\ &= E_x[E_{t|x}[L_{t,h(x)}]] \\ &= E_x\left[\sum_{t=1}^K L_{t,h(x)} P(t|x)\right] \end{aligned}$$

Choose  $h(x)=j$  s.t.  $\underline{\quad}$  is minimized

# Cancer example – one final look!

- Example: classify medical images as ‘cancer’ or ‘normal’

		Decision	
		cancer	normal
Truth	cancer	0	1000
	normal	1	0

# Cancer example – one final look!

- Example: classify medical images as ‘cancer’ or ‘normal’

		Decision	
		cancer	normal
Truth	cancer	0	1000
	normal	1	0

$P(N|x)$  vs.  $1000 P(C|x)$   
 expected loss if  $h(x)=C$  - do - if  $h(x)=N$

# Inference and decision: three approaches for classification

- Generative model approach:

- (I) Model  $p(x, C_k) = p(x|C_k)p(C_k)$

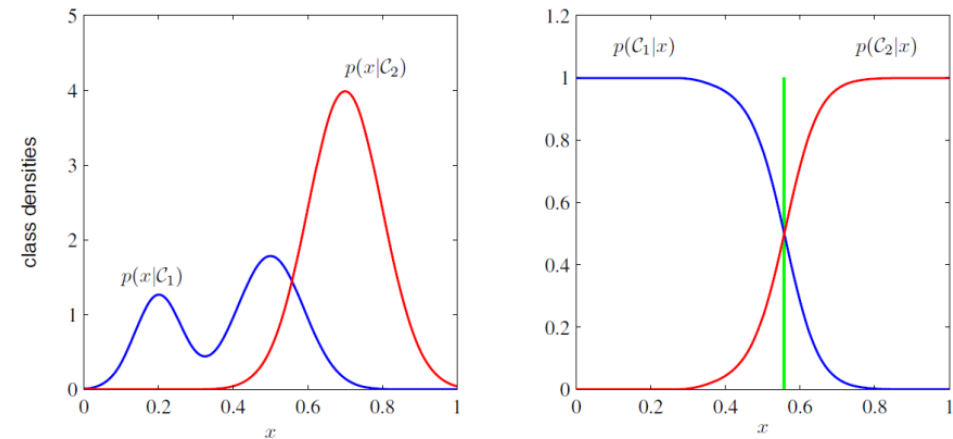
- (I) Use Bayes' theorem  $p(C_k|x) = \frac{p(x|C_k)p(C_k)}{p(x)}$

- (D) Apply optimal decision criteria

- Discriminative model approach:

- (I) Model  $p(C_k|x)$  directly

- (D) Apply optimal decision criteria



- Discriminant function approach:

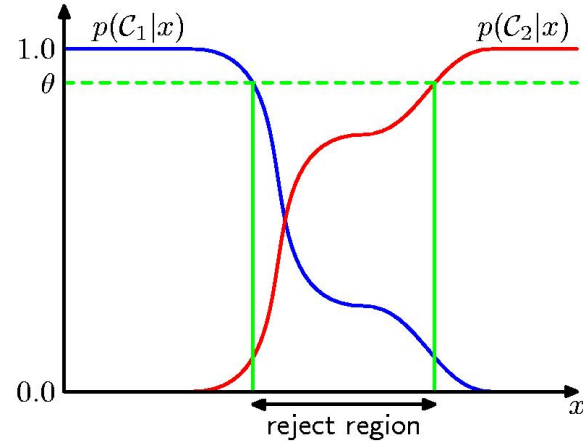
- (D) Learn a function that maps each  $x$  to a class label directly from training data

- Note: No posterior probabilities!

# Why separate Inference and Decision? (i.e., why infer (posterior) probabilities?)

- Minimizing risk (loss matrix may change over time)

- Reject option



- Combining models (Popular Naïve Bayes classifier)
- Etc.

# Problem Setting

- Naïve Bayes Classifier
  - Assumption: The features are independent given the class labels

# Problem Setting

- Naïve Bayes Classifier
  - Assumption: The features are independent given the class labels

Independent:  $p(X_1, X_2) = p(X_1)p(X_2)$

~~Conditionally~~ independent:  $p(X_1, X_2 | \underline{Y}) = p(X_1 | Y)p(X_2 | Y)$

# Naïve Bayes

Naive Bayes assumption:

$$\begin{aligned}p(X|Y) &= p(X_1, X_2, \dots, X_p | Y) \\&= p(X_p | X_1, X_2, \dots, X_{p-1}, Y) p(X_{p-1} | X_1, X_2, \dots, X_{p-2}, Y) \cdots p(X_1 | Y) \\&= p(X_p | Y) p(X_{p-1} | Y) \cdots p(X_1 | Y)\end{aligned}$$

$$\begin{aligned}p(Y|X) &= \frac{p(X_p | Y) p(X_{p-1} | Y) \cdots p(X_1 | Y) p(Y)}{p(X)} \\&\propto p(X_p | Y) p(X_{p-1} | Y) \cdots p(X_1 | Y) p(Y)\end{aligned}$$



# Naïve Bayes

- Assumption: The features are independent given the class labels
- Simple form for the probability distribution
- Not necessarily linear hyperplane 😊.
- Typically estimate by counting co-occurrences of feature value with class label
  - Maximum likelihood estimate
- Surprisingly powerful, especially in data with many features
  - High dimensional spaces

# Understanding Bayes Theorem

Given the data of accident reports and status as injured or not injured of the person after the accident.

$C_1 = \text{Inj}$   
 $C_2 = \text{NI}$

$$P(C_1) = \frac{6}{14}$$

$$P(C_2) = \frac{8}{14}$$

$$P(Y|C_1) = \frac{1}{2}$$

$$P(NR|C_1) = \frac{2}{3}$$

$$P(Y|C_2) = \frac{1}{8}$$

$$P(NR|C_2) = \frac{1}{4}$$

$$P(C_1 | Y, NR) \propto P(Y|C_1) P(NR|C_1) P(C_1) = \frac{1}{7}$$

$$P(C_2 | Y, NR) \propto P(Y|C_2) P(NR|C_2) P(C_2) = \frac{1}{5}$$

Case: Yamaha and Not repaired

Bike name	Repaired	Injured or Not injured
Yamaha	Yes	Injured
Yamaha	No	Injured
Suzuki	No	Not injured
TVS	Yes	Not injured
Honda	Yes	Not injured
Suzuki	Yes	Not injured
TVS	Yes	Injured
TVS	No	Injured
Honda	Yes	Not injured
Yamaha	No	Injured
Suzuki	Yes	Not injured
TVS	No	Injured
Honda	Yes	Not injured
Yamaha	No	Not injured

# Classification through Bayes Theorem

Given data on bikes and their features

$C = Y$   
 $C = H$

Bikes	weight	Engine
yamaha	100	300
yamaha	110	250
yamaha	92	250
yamaha	80	200
Honda	90	250
Honda	65	200
Honda	80	150
Honda	70	175

$p(85 | H)$   
 $p(w | H) \propto N(\mu, \sigma)$   
 $p(E | H)$

Predict the bike that was purchased from a given set of features,

Weight = 85 and engine = 250, Bike = ??

Where  $p(\text{yamaha}) = 0.5$  and  $p(\text{Honda}) = 0.5$

# Assumptions

- Weight and engine are continuous variables
- Weight and engine are independent variables

# Classification through Bayes Theorem

	Mean (weight)	Mean (Engine)	Variance (weights)	Variance (engine)
Yamaha(Y)	95.5	250	161	1666.66
Honda(H)	76.25	193.75	122.91	1822.91

Using Gaussian naïve Bayes,

$$P(Y/x(\text{weight}, \text{engine})) = p(Y) * p(\text{weight}/Y) * p(\text{engine}/Y) * (1/p(x))$$

$$P(H/x) = p(H) * p(\text{weight}/H) * p(\text{engine}/H) * (1/p(x))$$

Using Gaussian distribution,

Probability	weight	engine
Yamaha	0.022331	0.009775
Honda	0.026361	0.003924

$$p(x = v|c) = \frac{1}{\sqrt{2\pi\sigma_c^2}} e^{-\frac{(v-\mu_c)^2}{2\sigma_c^2}}$$

$$P(\text{yamaha}/x) > p(\text{Honda}/x)$$



**YAMAHA**