

M2. Decision Theory (incl. Bayes classifiers)

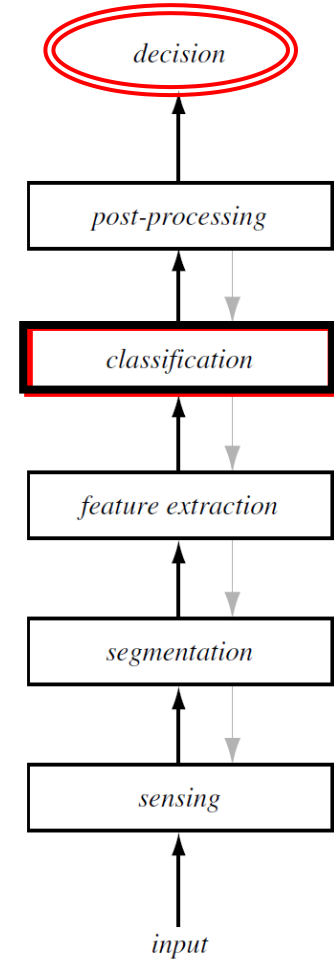
Manikandan Narayanan

Week 2 (Aug 4-)

PRML Jul-Nov 2025 (Grads Section)

Recall: Full PRML pipeline

- Before we delve into the ML parts, let's also look at decision/action, the final step!!
- Expected learning outcomes of this topic:
 - **primary:** Understand Decision Theory
 - Optimal Bayes classifier
 - Optimal regressor
 - **secondary:** Understand certain paradigms/terms in ML:
 - Density estimation (discriminative vs. generative modelling) in the context of supervised learning (classification/regression), and
 - use it to set the stage for unsupervised learning topics like clustering and supervised topics like *Naïve* Bayes classifier!



Acknowledgment of Sources

- Slides based on content from related
 - Courses:
 - IITM – Profs. Arun/Harish[HR]/Chandra[CC]/Prashanth’s PRML offerings (slides, quizzes, notes, etc.), Prof. Ravi’s “Intro to ML” slides – cited (e.g., [HR]/[HG]) 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]
 - Foundations of ML by Mohri, Rostamizadeh, and Talwalkar (content, figures, slides by Mohri, etc.). – [MRT]

Outline of Module M2

- M2. Decision Theory (incl. Bayes classifiers)
 - **M2.0 Decision Theory for Classification/Regression (common defns./notations)**
 - M2.1 Decision Theory for Classification (Bayes classifiers)
 - M2.2 Decision Theory for Regression (Squared loss, etc.)

M2.0 Decision Theory (for classification/regression)

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

- Inference step
 - Determine either $p(t|x)$ or $p(x, t)$. (density estimation)
- Decision step
 - For any given x , determine optimal t .
 - Optimality wrt (*empirical*) *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
 - $h: \mathcal{X} \rightarrow \{C_1, \dots, C_K\}$
 - $f: \mathcal{X} \rightarrow \mathbb{R}$

Handwritten red notes illustrating the feature vector \mathbf{x} and its components. The notes show $\mathbf{x} = (x_1, x_2, \dots, x_D)$ and a function $f(x)$ applied to each component, resulting in $f(x_1), f(x_2), \dots, f(x_D)$.

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 , $P(x) = P(t = C_1)P(x|C_1) + p(t = C_2)P(x|C_2)$
 $= P(C_1)P(x|C_1) + P(C_2)P(x|C_2)$

Outline of Module M2

- M2. Decision Theory (incl. Bayes classifiers)
 - M2.0 Decision Theory for Classification/Regression (common defns./notations)
 - **M2.1 Decision Theory for Classification (Bayes classifiers)**
 - M2.2 Decision Theory for Regression (Squared loss, etc.)

M2.1 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)

STOP & THINK: What is your guess for the optimal classifier (for binary classification)?

- That is, you are given a particular datapoint x .
- You already know $p(t = C_1 | x)$ and $p(t = C_2 | x)$ (say 0.3 and 0.7 respectively). Using this information,
 - how will you decide the optimal class label t for x ?
 - Will your prediction be $h(x) = C_1$ or $h(x) = C_2$?
 - What is the expected 0-1 loss for each of these two cases?

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 proof – minimum misclassification rate -- in equations

$$E[L] = P(\text{error}) = P(h(x) \neq t)$$

$$= \int \sum_{t=c_1, c_2} P(x, t) \mathbb{1}_{\{h(x) \neq t\}} dx$$

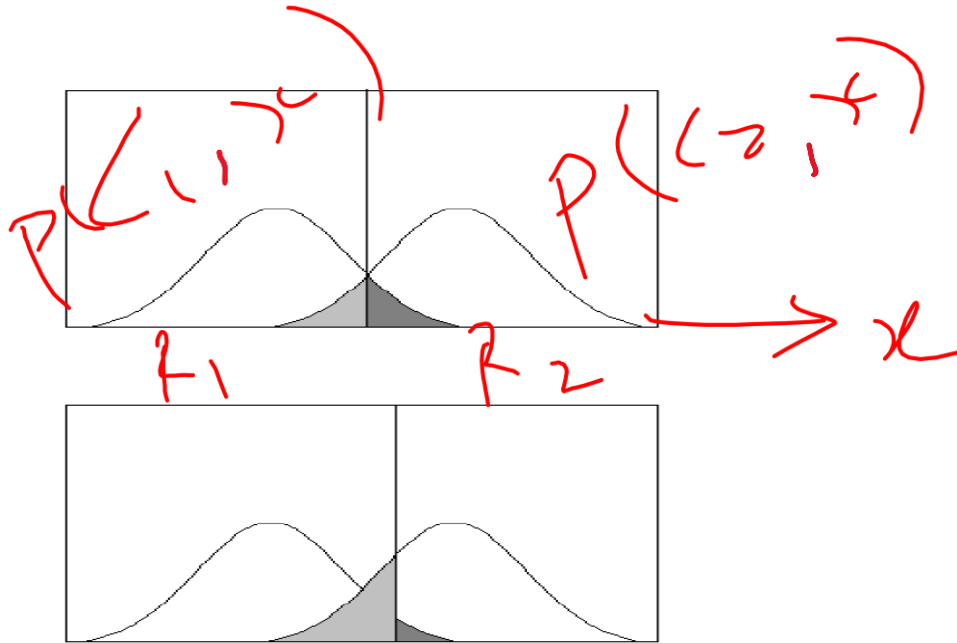
$$= \int \sum_t P(t|x) \mathbb{1}_{\{h(x) \neq t\}} P(x) dx$$

<min. by Bayes classifier
in prev. slide>

Optimality proof – minimum misclassification rate – in pictures

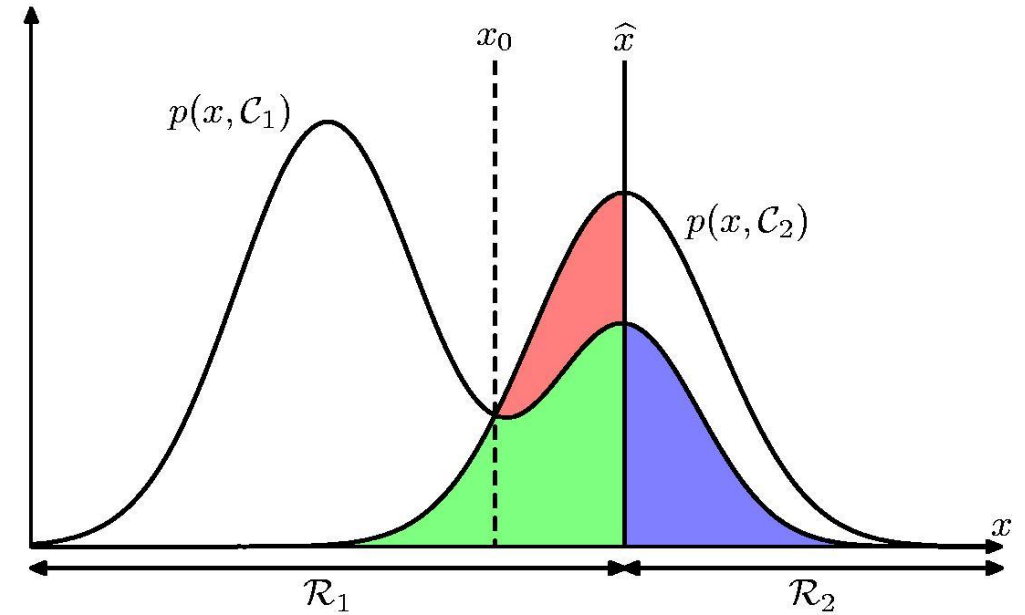
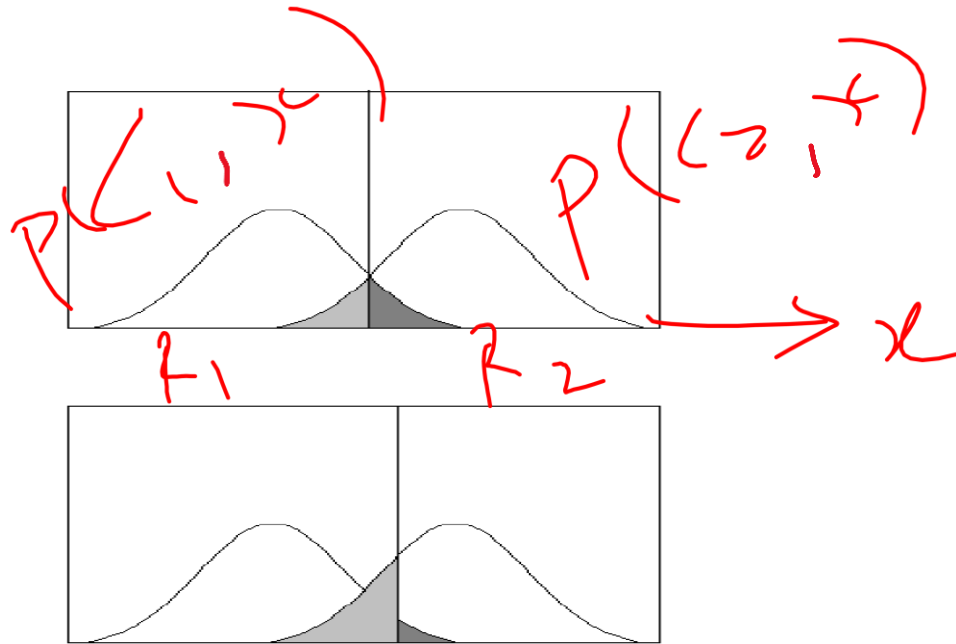
- **Goal:** Find $h: \mathcal{X} \rightarrow \{C_1, \dots, C_K\}$ s.t. $R(h) = E[L(h)]$ is minimized.
- That is, find optimal classifier $h^* = \arg \min_h R(h)$
- Let Decision region $R_i := \{x \in \mathcal{X} \mid h(x) = C_i\}$

Optimality - minimum misclassification rate



small note: $\mathbf{P}(t, \mathbf{x}) = P(t | \mathbf{x}) P(\mathbf{x}) \propto \mathbf{P}(t|\mathbf{x})$

Optimality - minimum misclassification rate



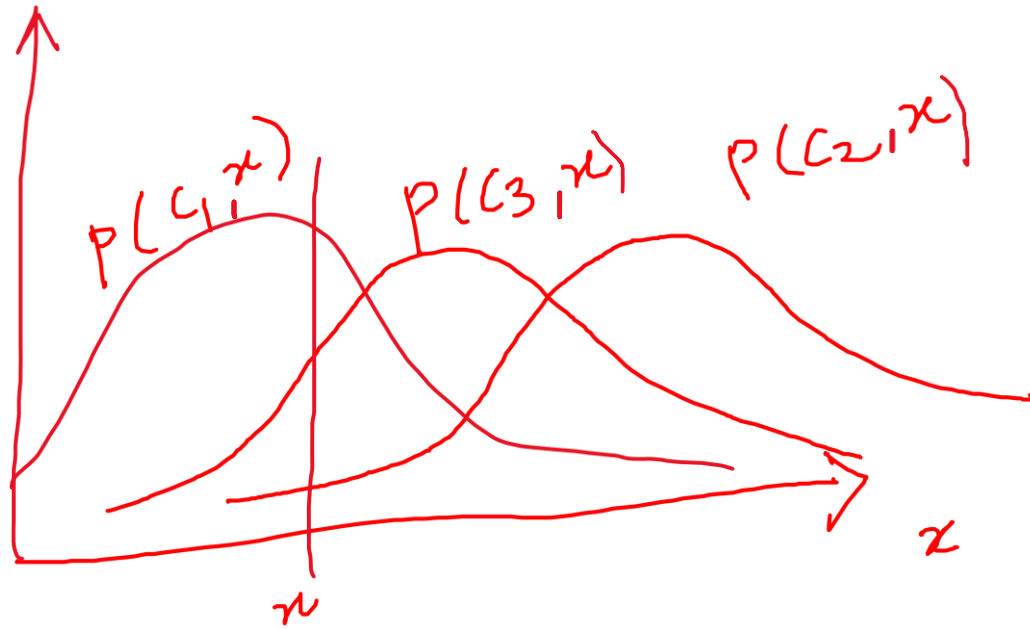
small note: $\mathbf{P}(\mathbf{t}, \mathbf{x}) = P(t | \mathbf{x}) P(\mathbf{x}) \propto \mathbf{P}(\mathbf{t} | \mathbf{x})$

Can decision regions be discontinuous in the optimal classifier?

Can decision regions be discontinuous in the optimal classifier?



What about $K > 2$ classes?



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

- $h(\mathbf{x}) = C_j$ if $P(t = C_j | \mathbf{x}) \geq P(t = C_{j'} | \mathbf{x}) \quad \forall j' \in \{1, \dots, K\} \setminus \{j\}$
 $= \operatorname{argmax}_{C_j} P(t = C_j | \mathbf{x})$ (ties broken arbitrarily)
- 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 – minimum misclassification rate --
in equations (for $K > 2$ classes)

$$E[L] = P(\text{error}) = P(h(x) \neq t)$$

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

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

<min. by Bayes classifier
in prev. slide>

Stop and Think! What have we seen so far?

- Optimizing mis-classification rate in $K=2$ and $K > 2$ settings.
- How about optimizing expected loss for general loss functions?
 - From 0—1 loss matrix to general loss matrix!

Bayes classifier – General Loss Function (Matrix)

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

		Decision $h(x)$	
		cancer	normal
Truth	cancer	0	1000
	normal	1	0

Bayes classifier – General Loss Function (Matrix)

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

		Decision $h(x)$	
		cancer	normal
Truth	cancer	0	1000
	normal	1	0

Intuition:

$$h(x) = \begin{cases} N(\text{normal}) & \text{if } P(N|x) > (0.5)P(C|x) \\ C(\text{cancer}) & \text{o.w.} \end{cases}$$

Bayes classifier – General Loss Function (Matrix)

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

\times

		Decision $h(x)$	
		cancer	normal
Truth	cancer	0	1000
	normal	1	0

Intuition:

$$h(x) = \begin{cases} N(\text{normal}) & \text{if } P(N|x) > 1000 P(C|x) \\ C(\text{cancer}) & \text{o.w.} \end{cases}$$

$L = \begin{matrix} & \begin{matrix} 1 & 2 & \dots & K \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ \vdots \\ K \end{matrix} & \begin{bmatrix} 0 & 0 & \dots & 1000 \\ 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ L_{ij} & \dots & \dots & 0 \end{bmatrix} \end{matrix}$

(Note: L_{ij} is circled in the matrix)

Optimality - Minimum Expected Loss (integration notation)

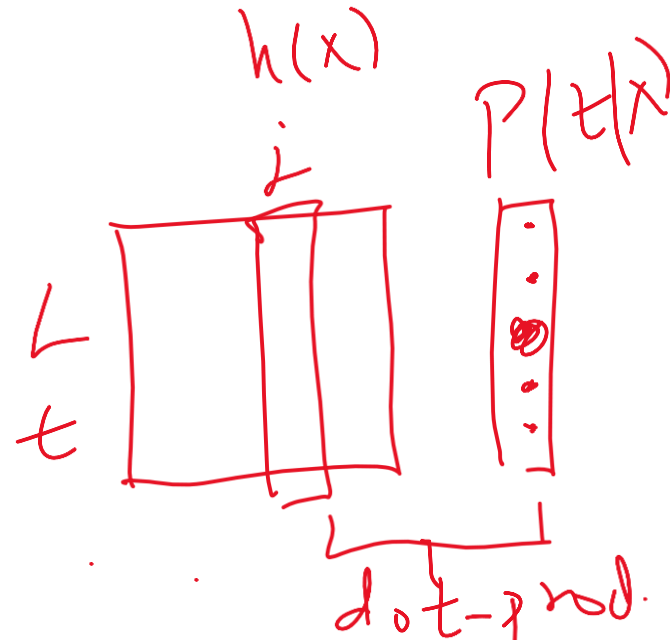
$$\begin{aligned} \rightarrow E_{x,t}[L] &= \int \sum_t L_{t,h(x)} P(x,t) dx \\ &= \int \underbrace{\sum_{t \in C} P(t|x)}_{\substack{\text{Choose} \\ h(x)=j \text{ s.t. } \searrow \text{ is minimized}}} L_{t,h(x)} P(x) dx \\ &\quad E_{t|x}[L_{t,h(x)}] \end{aligned}$$

Optimality - Minimum Expected Loss (contd. expectation notation)

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

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


Optimality - Minimum Expected Loss – Minimize $E_{t|x}[L_{t,j}]$ (aka dot-prod.) over all j .



$$h(x) = \arg \min_j \sum_t L_{t,j} \cdot P(t|x)$$

Cancer example – one final look!

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

		Decision		
		cancer	normal	
Truth	cancer	0	1000	$P(t x)$ <div style="border: 1px solid red; padding: 2px; display: inline-block;">  </div>
	normal	1	0	

$$\underbrace{P(N|x)}_{\text{exp. loss}} \text{ vs. } \underbrace{1000 P(C|x)}_{\text{exp. loss}}$$

if $h(x)=C$ if $h(x)=N$

matches intuition:

$$h(x) = \begin{cases} N(\text{normal}) & \text{if } P(N|x) > 1000 P(C|x) \\ C(\text{cancer}) & \text{o.w.} \end{cases}$$

Inference and decision: three approaches for classification

- Generative model approach:

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

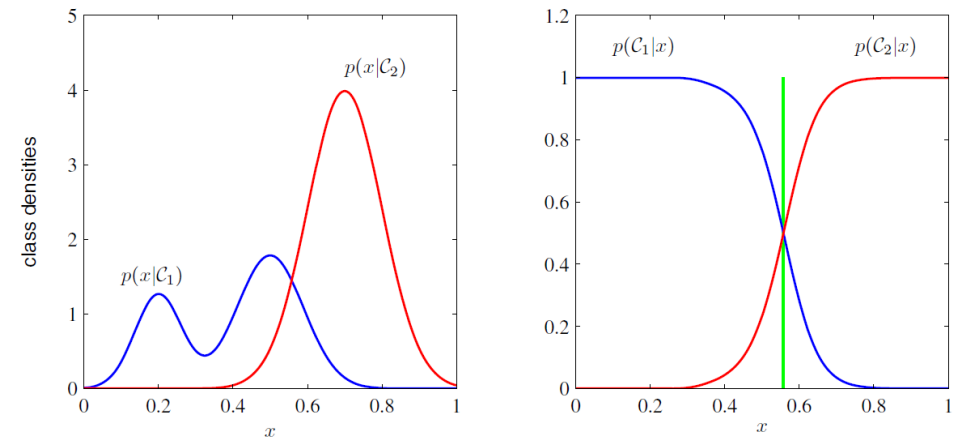
- (I) Use Bayes' theorem $p(C_k|x) = \frac{p(x|C_k)p(C_k)}{p(x)}$

- (D) Apply optimal decision criteria

- Discriminative model approach:

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

- (D) Apply optimal decision criteria



- Discriminant function approach:

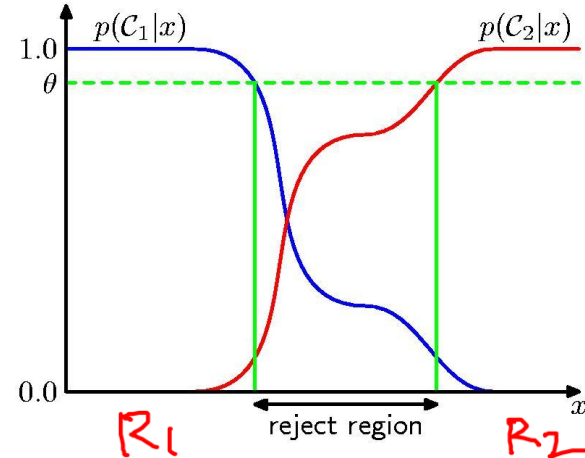
- (D) Learn a function that maps each x to a class label directly from training data

- 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.

Example of generative vs. discriminative models for the same task

- Task: Classification
 - Discriminative model $p(t|x)$: Logistic regression
 - Generative model $p(t,x) = p(t) p(x|t)$: Naïve Bayes classifier
- In general, discriminative model preferred – folklore
 - But a nuance: discriminative preferred for large sample sizes, vs. generative for smaller sample sizes (IF model assumptions are satisfied!)
[Optional reference: On Discriminative vs. Generative Classifiers: A comparison of logistic regression and naive Bayes. *Andrew Ng and Michael Jordan. NIPS 2001.*]

Final notes for tutorial

- **Exercise:** What is Precision, Recall, etc., in terms of Probab. over (x, t) ?
- **Loss matrix \neq Confusion matrix,**
 - but to estimate $R[h]=E[L(h)]$ of a learned (trained) classifier $h(\cdot)$, we can use both matrices.
 - **Important:** To estimate $R[h]$, use confusion matrix wrt test data. Why??

Some examples

- See Bayes classifier examples in [HG]Notes!

Outline of Module M2

- M2. Decision Theory (incl. Bayes classifiers)
 - M2.0 Decision Theory for Classification/Regression (common defs./notations)
 - M2.1 Decision Theory for Classification (Bayes classifiers)
 - **M2.2 Decision Theory for Regression (Squared loss, etc.)**

M2.2 Decision Theory for Regression

- Inference step
 - Determine $p(\mathbf{x}, t)$ or $p(t \mid \mathbf{x})$.
- Decision step
 - For given \mathbf{x} , make optimal prediction $f(\mathbf{x})$ for t (min. risk or expected loss).
 - Given a loss fn., $E[L] = \iint L(t, f(\mathbf{x})) p(\mathbf{x}, t) d\mathbf{x} dt$

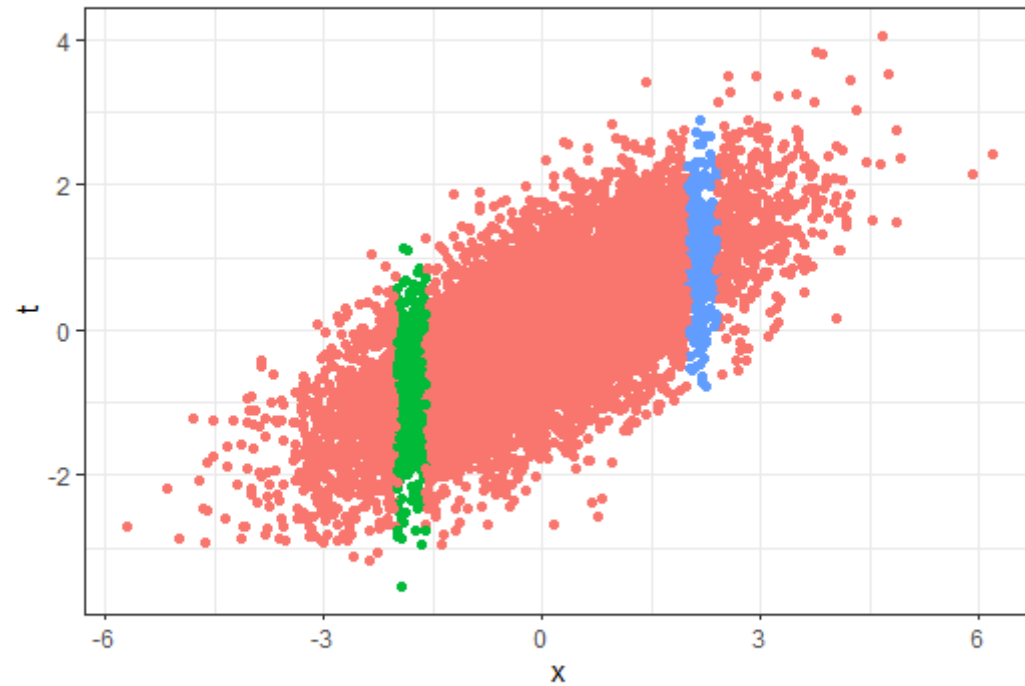
The Squared Loss Function

$$E[L] = \iint (f(\mathbf{x}) - t)^2 p(\mathbf{x}, t) d\mathbf{x} dt$$

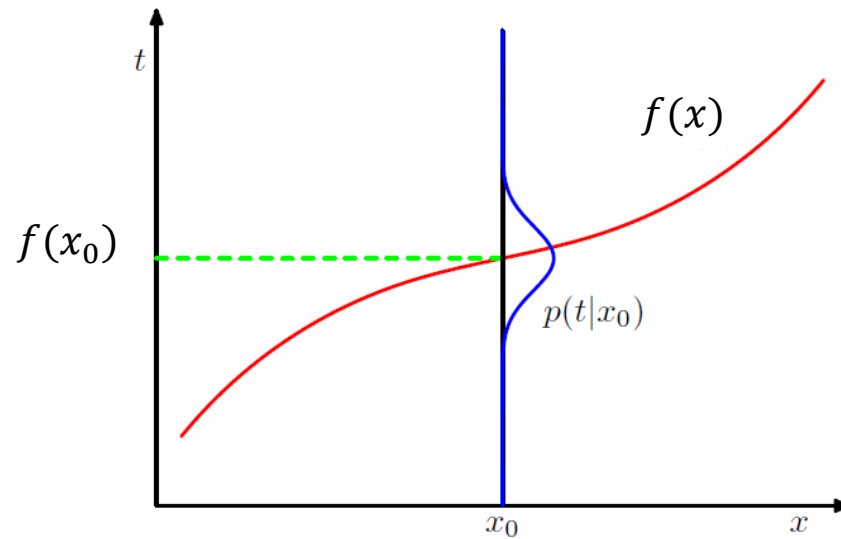
$$E[L] = \int \int_{\mathbf{x}} (f(\mathbf{x}) - t)^2 p(t|\mathbf{x}) dt p(\mathbf{x}) d\mathbf{x}$$

$$E_{\mathbf{x}, t}[L] = E_{\mathbf{x}} \left[E_{t|\mathbf{x}} [(f(\mathbf{x}) - t)^2] \right]$$

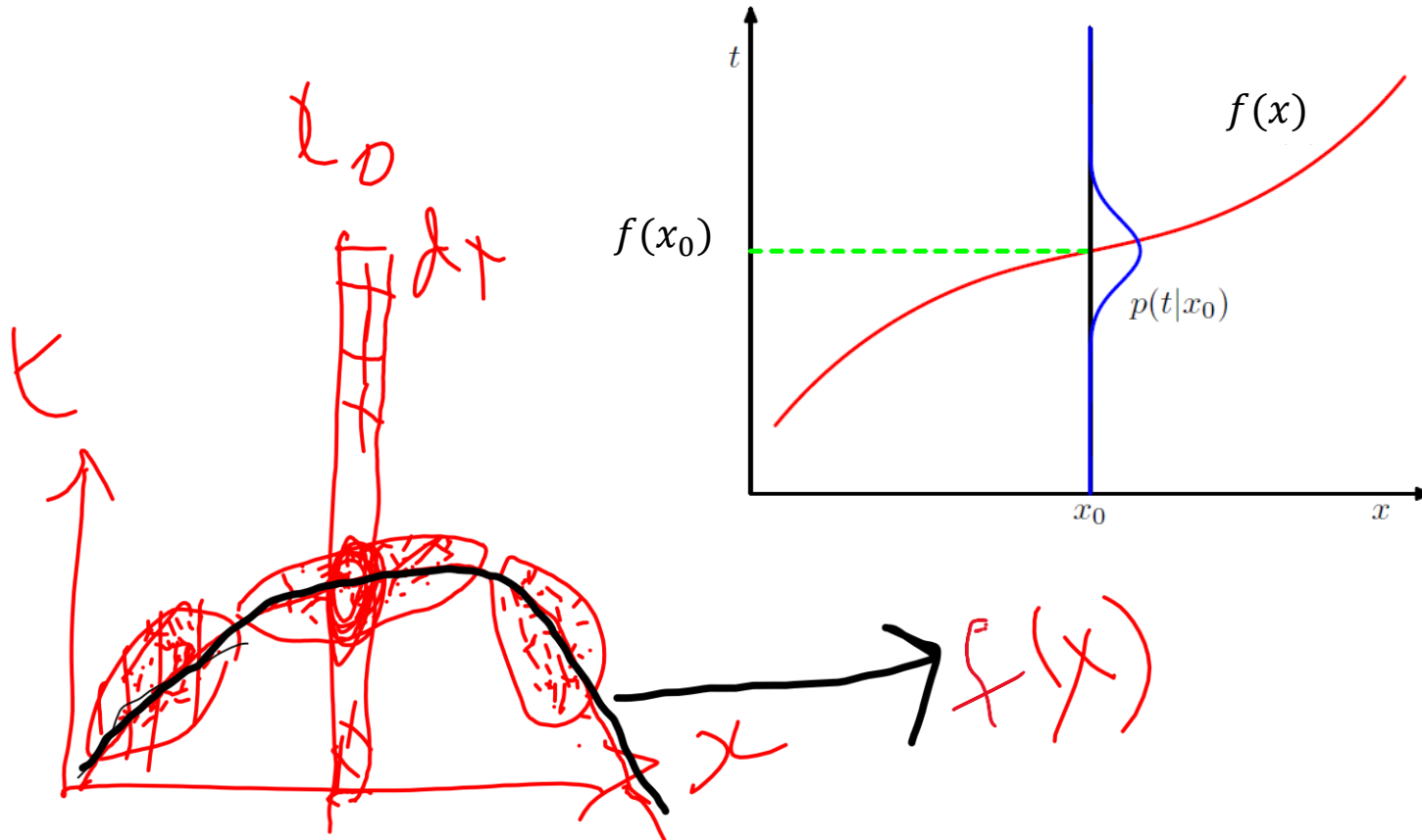
What would be a good minimizer if (x,t) is bivariate normally distributed like in our previous example $(\mathbf{X} = W_1 + W_2, t = W_2)$? (with w_1, w_2 indept. std. Gaussian rvs.)



What if (x,t) has some general distribution?



What if (x,t) has some general distribution?



$$\int p(t|x) t dt$$

$E[t|x]$

Optimality – min. squared loss (conditional expectation as a minimizer)

$$E[L] = E_x \left[E_{t|x} \left[(f(x) - t)^2 \right] \right]$$

$$\{f(x) - t\}^2 = \{f(x) - E[t|x] + E[t|x] - t\}^2$$

$$= \{f(x) - E[t|x]\}^2 + 2\{f(x) - E[t|x]\}\{E[t|x] - t\} + \{E[t|x] - t\}^2$$

$$E_{t|x}[(f(x) - t)^2] =$$

Optimality – min. squared loss (conditional expectation as a minimizer)

$$E[L] = E_x \left[E_{t|x} \left[(f(x) - t)^2 \right] \right]$$

$$\{f(x) - t\}^2 = \{f(x) - E[t|x] + E[t|x] - t\}^2$$

$$= \underbrace{\{f(x) - E[t|x]\}^2}_{\text{signal}} + \underbrace{2\{f(x) - E[t|x]\}\{E[t|x] - t\}}_{\text{noise}} + \underbrace{\{E[t|x] - t\}^2}_{\text{noise}}$$

$$E_{t|x} \left[(f(x) - t)^2 \right] = \text{signal} + \text{noise}$$

$$\boxed{\text{Var}(t|x)}$$

$$E[L] = \int \{f(x) - E[t|x]\}^2 p(x) dx + \underbrace{\int \text{var}(t|x) p(x) dx}_{\text{noise}}$$

$$y(x) = E[t|x]$$

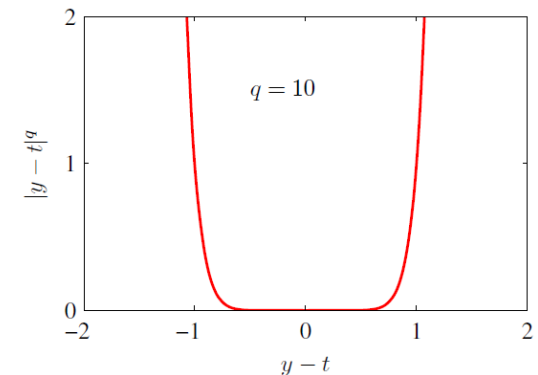
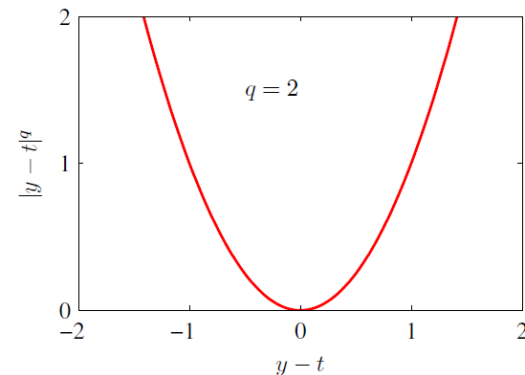
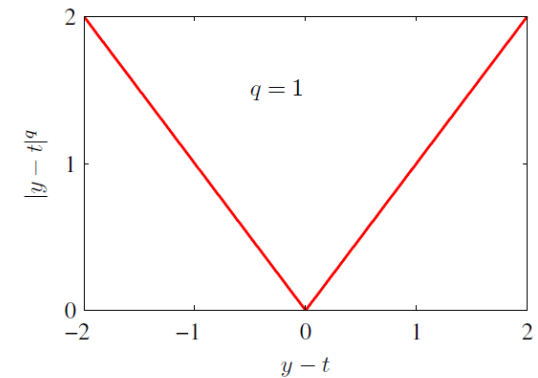
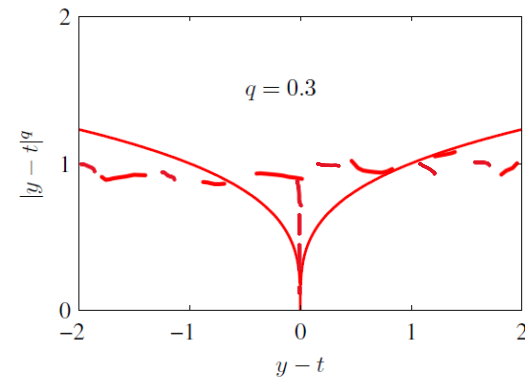
noise

From Squared to Minkowski loss

$$E[L_q] = \iint |f(\mathbf{x}) - t|^q p(\mathbf{x}, t) d\mathbf{x} dt$$

Conditional mean / median / mode
for $q = 2$ / $q=1$ / $q \rightarrow 0$
respec.

$p(t|\mathbf{x})$ – inherent variab. in data



Three approaches again (for regression)

- Generative model approach:
 - (I) Model $p(t, \mathbf{x}) = p(\mathbf{x}|t)p(t)$
 - (I) Use Bayes' theorem $p(t|\mathbf{x}) = \frac{p(\mathbf{x}|t)p(t)}{p(\mathbf{x})}$
 - (D) Take conditional mean/median/mode/any other optimal decision outcome as $f(\mathbf{x})$
- Discriminative model approach:
 - (I) Model $p(t|\mathbf{x})$ directly
 - (D) Take conditional mean/median/mode/any other optimal decision outcome as $f(\mathbf{x})$
- Direct regression approach:
 - (D) Learn a regression function $f(\mathbf{x})$ directly from training data

In summary

- Decision theory (for classifn./regn.):
 - inference and decision steps
 - Generative vs. discriminative models for inference
 - Minimum risk (expected loss) for optimal decision
 - Different loss functions possible for classification; similarly for regression; but require knowledge of posterior (or joint directly or via prior and liklhd.) densities
 - Direct/discriminant approach also possible to take decision
- Optimal models, but how to learn them?
 - Bayes classifier for classifn. and conditional mean for regn., PROVIDED we've access to the joint or posterior distbn.
 - Next → Q: How do we learn these distbns. (and hence the classifier/regressor) from data? A: Density estimation of $P(X,Y)$ or $P(Y|X)$.

Thank you!

Backup slides follow

Backup

Optimality of multi-class classifier (max. accuracy)

$$P(\text{error}) = P_{x,t}(h(x) \neq t) = 1 - \overbrace{P_{x,t}(h(x) = t)}^{P(\text{correct})}$$

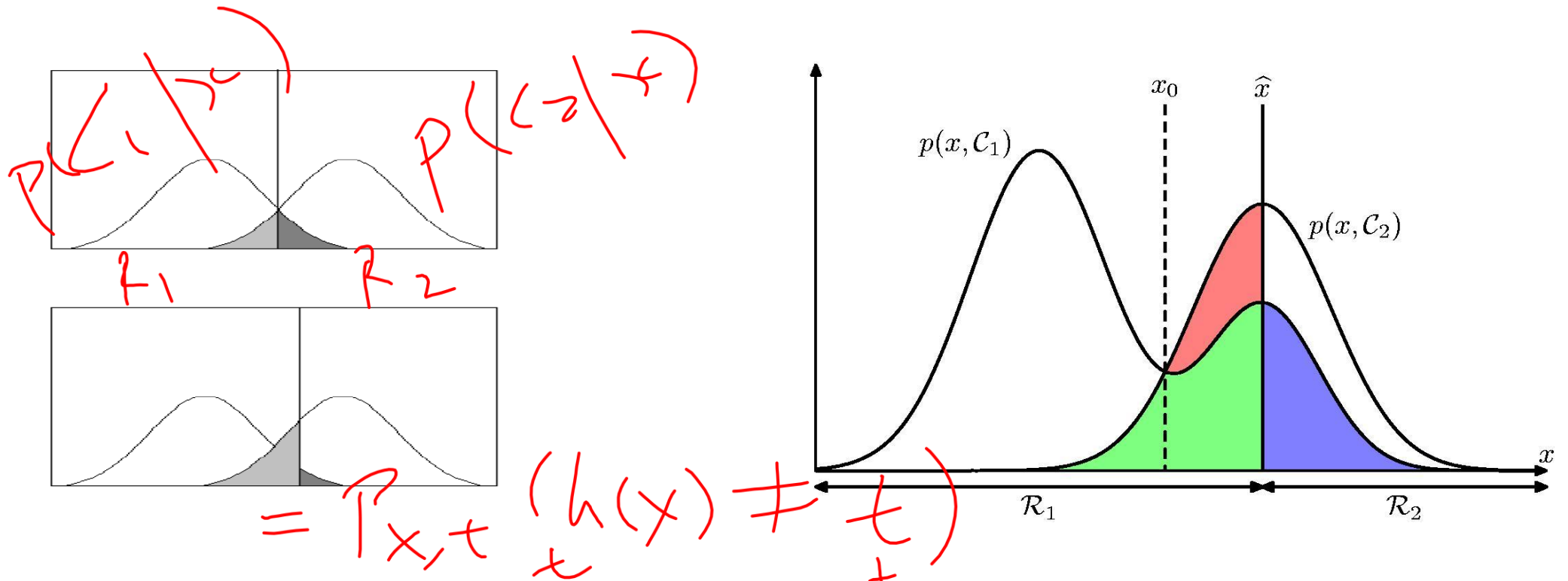
$$\mathbb{1}_C = \begin{cases} 1 & \text{if } C \text{ is } T \\ 0 & \text{if } C \text{ is } F \\ & \text{else} \end{cases}$$

$$P(\text{correct})$$

$$\int_{\mathcal{X}} \sum_{t=1}^C P(x, t) \mathbb{1}_{\{h(x)=t\}} dx = \int_{\mathcal{X}} \left(\sum_{t=1}^C P(t|x) \mathbb{1}_{\{h(x)=t\}} \right) P(x) dx$$

[CMB]

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 of multi-class classifier (min. error)

$$P(\text{error}) = \left(\int \sum_{t=1}^K p(x, t) \mathbb{1}_{(h(x) \neq t)} dx \right)$$

$$p(x, t=1) + \underbrace{p(x, t=2)}_{[P(x) - \underline{p(x, t=2)}]} + \dots + p(x, t=K)$$

Optimality - Minimum Expected Loss (indicator fn. notation)

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

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