# Foundations of Machine Learning
# AI2000 and AI5000

FoML-11
Bayesian Regression

FoML-12

Decision Theory

Dr. Konda Reddy Mopuri
Department of AI, IIT Hyderabad
July-Nov 2025

भारतीय सांकेतिक विज्ञान संस्थ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad

DiL
Data-driven Intelligence
& Learning Lab

# So far in FoML

- What is ML and the learning paradigms

- Probability refresher

- MLE, MAP, and fully Bayesian treatment

- Linear Regression with basis functions - and regularization

- Model selection

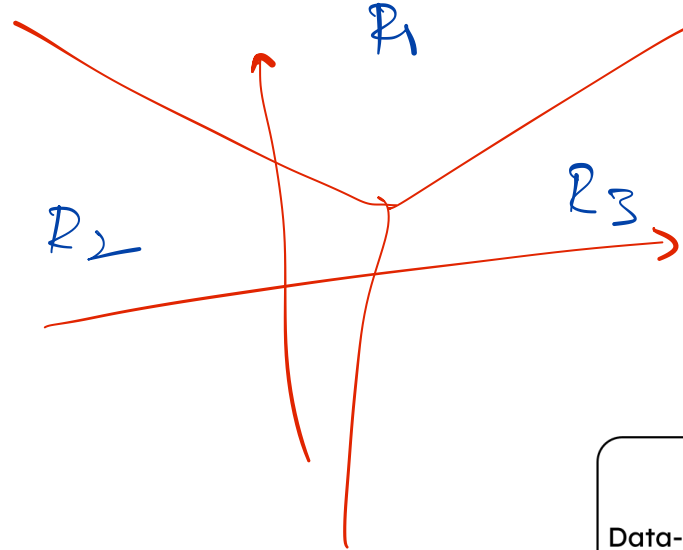- Bias-Variance Decomposition/Tradeoff (Bayesian Regression)

# Decision Theory

# Decision Theory

- Dataset: i/p vectors $\mathbf{x} \in \mathbb{R}^D$, ground truth $t \in \{C_1, C_2, \ldots, C_K\}$
- Divide the i/p space $\mathbb{R}^D$ into K decision regions $R_k$, $k = \{1, 2, \ldots K\}$
- For every data point
  - Ground truth $t_n$
  - Prediction $y(x_n, w) = \hat{t}_n$

$R_1$

$R_2$

$R_3$

# Decision Theory

- Confusion Matrix



Diagonal elements - correct predictions

Off-diagonal elements - incorrect predictions

# Decision Theory - Misclassification Rate

- Goal of classification - Minimize the misclassification rate

- Assume the data are drawn independently from the joint distribution

- Probability of a misclassification: $p(\text{mistake}) = \sum_{i=1}^{K} \sum_{k \neq i} p(\mathbf{x} \in R_i, C_k)$

$$p(\text{mistake}) = 1 - P(\text{Correct classification})$$

$$= 1 - \sum_{k=1}^{K} P(x \in R_k, C_k)$$

# Decision Theory - Misclassification Rate

- Minimizing the misclassification rate    (how to ensure this?)

  - Assign x to class $C_k$ if    $p(\mathbf{x}, t = C_k) > p(\mathbf{x}, t = C_j), \ \ \forall j \neq k$

  - We know that

$$P(\underline{x}, C_k) = P(C_k | \underline{x}) \cdot P(\underline{x})$$

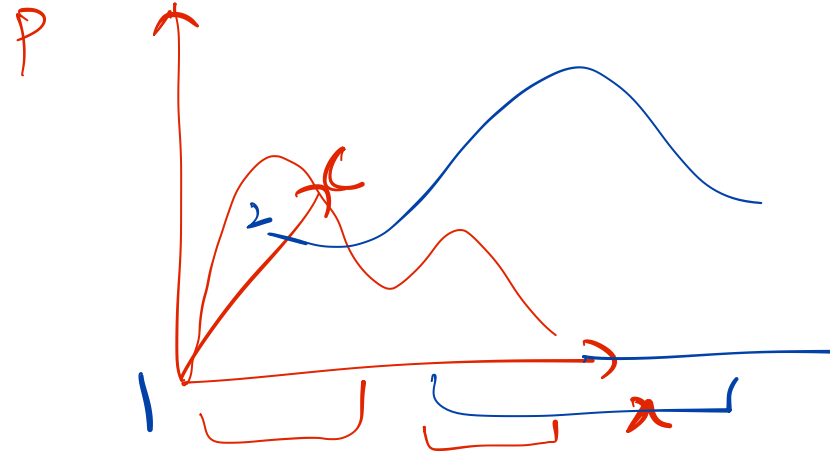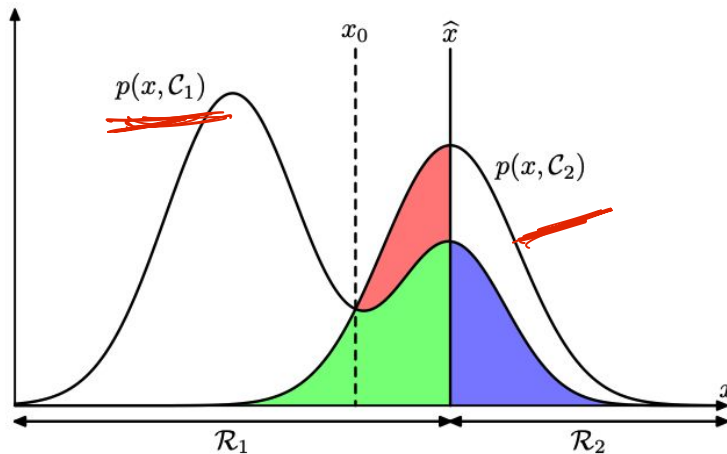look for the largest
posterior prob $P(C_k | \underline{x})$

$P(x, C_j)$

$j \in \{1 \ldots k\}$

# Decision Theory - Misclassification Rate

# Minimizing the Misclassification Rate - Issues

- Not all errors have the same impact!

- E.g. medical diagnosis
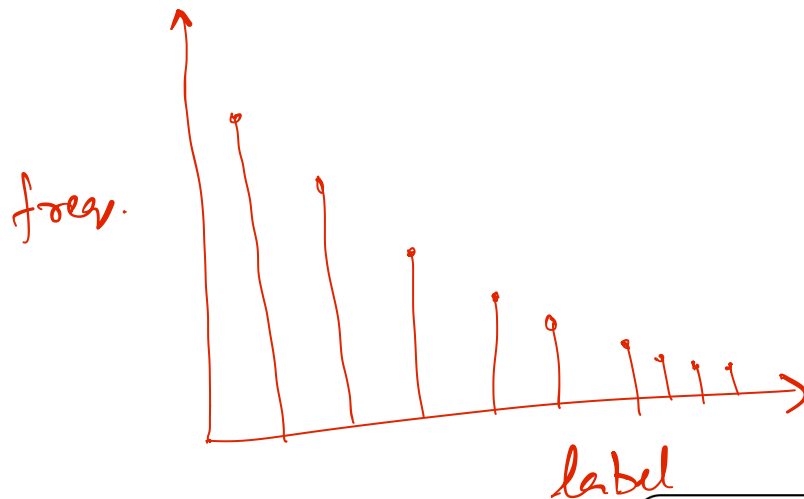
  - E1:

  - E2: ✓

$$P(\underline{x}, \hat{t} = D) > P(\underline{x}, \hat{t} = H) \quad \text{when } t = H$$

$$P(\underline{x}, \hat{t} = H) > P(\underline{x}, \hat{t} = D) \quad \text{when } t = D$$

# Minimizing the Misclassification Rate - Issues

- ## Class imbalance
    - ○ May lead to skewed view of the classifier's performance

1 %

freq.

label

# Expected Loss

- Possible solution: use different weights for different error types

$$\text{pred} \longrightarrow$$

$$
\mathbf{L} = \begin{pmatrix} 0 & 1000 \\ 1 & 0 \end{pmatrix} \quad \begin{matrix} C \\ H \end{matrix} \quad \Big\downarrow \text{G.T.}
$$

$$\sum_{k} \sum_{j} \int_{R_j} L_{k,j} \, p(x, C_k) \, dx$$

another way to write

$$\mathbb{E}[L] = \sum_{k,j} L_{k,j} \int_{\mathcal{R}_j} p(x, C_k) \, dx$$

Minimize the expected loss: (assign x to Ck if )   $\sum_{j=1}^{K} L_{k,j} \, P(C_k | x)$   is minimal

# Classification Strategies

- Discriminant functions
  - Direct functions of i/p to target $\quad t = y(\mathbf{x}, \mathbf{w})$
- Probabilistic Discriminant models
  - Posterior class probabilities $\quad p(C_k / \mathbf{x})$
- Probabilistic generative models
  - Class-conditional models $\quad p(\mathbf{x}/C_k)$
  - Prior class probabilities $\quad p(C_k)$

# Next
# Probabilistic Generative Models