# Foundations of Machine Learning AI2000 and AI5000

FoML-16
Least Squares for Regression

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

భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad

DiL
Data-driven Intelligence
& Learning Lab

# So far in FoML

- Intro to ML and Probability refresher

- MLE, MAP, and fully Bayesian treatment

- Supervised learning

    a. Linear Regression with basis functions (regularization, model selection)

    b. Bias-Variance Decomposition (Bayesian Regression)

    c. Decision Theory - three broad classification strategies

        ■ Probabilistic Generative Models - Continuous & discrete data

        ■ Discriminant Functions

భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
**Indian Institute of Technology Hyderabad**

DiL
**Data-driven Intelligence**
**& Learning Lab**

# Least Squares for Classification

# Least Squares for Classification

- Consider K classes
- Each class 'k' has its own linear model $y_k(\mathbf{x}) = w_k^T \mathbf{x} + w_{k0}$

# Least Squares for Classification

$$T = \begin{bmatrix} - t_1^T - \\ \vdots \\ - t_N^T - \end{bmatrix}_{N \times K}$$

- Shorter notation $y(\mathbf{x}) = \widetilde{\mathbf{W}}^T \tilde{\mathbf{x}}$

$$t_\eta = (0\,0 \cdots 1\,0\,0)^T_{1 \times K}$$

$$\widetilde{\mathbf{W}} = \begin{bmatrix} w_{10} - w_1 - \\ \vdots \\ w_{K0} - w_K - \end{bmatrix}_{K \times M}^T$$

Assign x to $C_k$, where

$$\tilde{\mathbf{x}} = (1, \underline{x})^T_{1 \times M} \qquad X = \begin{bmatrix} \tilde{x}_1 \\ \vdots \\ \tilde{x}_N \end{bmatrix}_{N \times M} \qquad k = \underset{j}{\arg\max}\; y_j(\underline{x})$$

$$y(\mathbf{x}) = \begin{bmatrix} y_1(\underline{x}) \\ \vdots \\ y_K(\underline{x}) \end{bmatrix}_{K \times 1}$$

on single sample

$$E(\underline{x}_\eta) = \left\| y(\underline{x}_\eta) - t_\eta \right\|^2$$

# Least Squares for Classification

- Data matrix $X_{N \times M}$ $\begin{bmatrix} \text{row is a} \\ \text{sample} \end{bmatrix}$

- Target matrix $T_{N \times K}$

Parameter matrix $\widetilde{W}_{M \times K}$ $\begin{bmatrix} \text{column is} \\ \text{per class discriminant} \end{bmatrix}$

Use regression (sum of squares) error function

$$E_D(\widetilde{\mathbf{W}}) = \frac{1}{2} \sum_{n=1}^{N} \sum_{k=1}^{K} \sum_{m=1}^{M} \left( t_{nk} - x_{nm} \widetilde{w}_{mk} \right)^2$$

భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad

DiL
Data-driven Intelligence
& Learning Lab

# Least Squares for Classification

The error function can be conveniently written as

$$E_D(\widetilde{\mathbf{W}}) = \frac{1}{2}\text{Tr}\left\{(\widetilde{\mathbf{X}}\widetilde{\mathbf{W}} - \mathbf{T})^{\text{T}}(\widetilde{\mathbf{X}}\widetilde{\mathbf{W}} - \mathbf{T})\right\}$$

Minimize $E_D(\widetilde{\mathbf{W}})$ as a function of $\widetilde{\mathbf{w}}$:

$$(\widetilde{X}^T\widetilde{X})^{-1}\widetilde{X}^T \cdot T$$

$$P(c_k|x)$$

$$\frac{\partial}{\partial \widetilde{w}}(\cdot) = 0$$

the familiar pseudo inverse formulation

# Least Squares Issues - Outliers



Magenta - LS classifier
Green - Logistic Regression classifier

# Least Squares Issues - Masking

$$\begin{bmatrix} \cdot & - & \frac{1}{\ast} & \cdot \end{bmatrix}$$

prediction for the correct class is expected to be close to 1

$$\frac{y|x\rangle}{\|w\|}$$



Left - LS classifier
Right - Logistic Regression classifier

భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad

DiL
Data-driven Intelligence
& Learning Lab

# Least Squares Issues - Predictions ≠ Probabilities

$\mathbf{y}_{LS}(\mathbf{x})$   are not probabilities

If   $\underset{\sim}{a}^T \underline{t_n} + b = 0$

$\Rightarrow \quad a^T y(\underline{x_n}) + b_n = 0$

with

$\underline{a} = [1\ 1111]$

$b = -1$

we know that

$\underset{\sim}{a}^T \underline{t_n} + b = 0$

holds ∀ $n$

this implies

$\sum_{i=1}^{K} y(\underline{x_n})_i = 1$ ∀ $n$

భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
**Indian Institute of Technology Hyderabad**

**DiL**
Data-driven Intelligence
& Learning Lab

Rough

# Next
# The Perceptron