# Fundamentals of Machine Learning

## MODEL FITTING & PARAMETER ESTIMATION

## LEAST SQUARES

Amit K Roy-Chowdhury

# Model Fitting / Training

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \; \mathcal{L}(\boldsymbol{\theta})$$

Loss function / objective function

# Maximum likelihood estimation

$$\hat{\boldsymbol{\theta}}_{\text{mle}} = \underset{\boldsymbol{\theta}}{\arg\max} \sum_{n=1}^{N} \log p(\boldsymbol{y}_n | \boldsymbol{x}_n, \boldsymbol{\theta})$$

Estimated
Parameter

Probability

Model

Since most optimization algorithms are designed to minimize cost functions, we redefine the objective function to be the (conditional) negative log likelihood or NLL and we minimize NLL

$$\text{NLL}(\boldsymbol{\theta}) \triangleq -\log p(\mathcal{D} | \boldsymbol{\theta}) = -\sum_{n=1}^{N} \log p(\boldsymbol{y}_n | \boldsymbol{x}_n, \boldsymbol{\theta})$$

# Notation and Form for MAP

Notation: $\hat{\theta}_{MAP}$    maximizes the posterior PDF

$$\hat{\theta}_{MAP} = \arg\max_{\theta} p(\theta \,|\, \mathbf{x})$$

Equivalent Form (via Bayes' Rule): $\boxed{\hat{\theta}_{MAP} = \arg\max_{\theta} [\,p(\mathbf{x}\,|\,\theta)\,p(\theta)\,]}$

Proof: Use $\quad p(\theta \,|\, \mathbf{x}) = \dfrac{p(\mathbf{x}\,|\,\theta)\,p(\theta)}{p(\mathbf{x})}$

$$\hat{\theta}_{MAP} = \arg\max_{\theta} \left[ \frac{p(\mathbf{x}\,|\,\theta)\,p(\theta)}{p(\mathbf{x})} \right] = \arg\max_{\theta} [\,p(\mathbf{x}\,|\,\theta)\,p(\theta)\,]$$

# Least Squares Approach

All the previous methods we've studied… required a underline{probabilistic} model for the data: underline{Needed the PDF} $p(\mathbf{x};\boldsymbol{\theta})$

For a Signal + Noise problem we needed:
Signal Model & Noise Model

**Least-Squares is underline{not} statistically based!!!**
**$\Rightarrow$ Do underline{NOT need} a PDF Model**

# Linear Least Squares

A <u>linear</u> least-squares problem is one where the parameter observation model is linear: $\mathbf{s} = \mathbf{H}\theta$ $\qquad$ $\mathbf{x} = \mathbf{H}\theta + \mathbf{e}$

$N \times 1$

$N \times p$ Known Matrix

$p \times 1$

$p = $ **Order** of the model

We must assume that **H** is full rank… otherwise there are multiple parameter vectors that will map to the same **s!!!**

# Linear Least Squares

For the linear model the LS cost is:
$$J(\boldsymbol{\theta}) = \sum_{n=0}^{N-1} \big(x[n] - s[n;\boldsymbol{\theta}]\big)^2$$

$$= (\mathbf{x} - \mathbf{H}\boldsymbol{\theta})^T (\mathbf{x} - \mathbf{H}\boldsymbol{\theta})$$

Now, to minimize, first expand:

$$J(\boldsymbol{\theta}) = \mathbf{x}^T\mathbf{x} - \mathbf{x}^T\mathbf{H}\boldsymbol{\theta} - \boldsymbol{\theta}^T\mathbf{H}^T\mathbf{x} + \boldsymbol{\theta}^T\mathbf{H}^T\mathbf{H}\boldsymbol{\theta}$$

$$= \mathbf{x}^T\mathbf{x} - 2\mathbf{x}^T\mathbf{H}\boldsymbol{\theta} + \boldsymbol{\theta}^T\mathbf{H}^T\mathbf{H}\boldsymbol{\theta}$$

*Scalar = scalar*$^T$ So…
$\boldsymbol{\theta}^T\mathbf{H}^T\mathbf{x} = (\boldsymbol{\theta}^T\mathbf{H}^T\mathbf{x})^T = \mathbf{x}^T\mathbf{H}\boldsymbol{\theta}$

Now setting $\dfrac{\partial J(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \mathbf{0}$ gives $-2\mathbf{H}^T\mathbf{x} + 2\mathbf{H}^T\mathbf{H}\hat{\boldsymbol{\theta}} = \mathbf{0}$

Called the
"LS Normal Equations"

$$\mathbf{H}^T\mathbf{H}\hat{\boldsymbol{\theta}} = \mathbf{H}^T\mathbf{x}$$

Because $\mathbf{H}$ is full rank we know that $\mathbf{H}^T\mathbf{H}$ is invertible:

$$\hat{\boldsymbol{\theta}}_{LS} = \left(\mathbf{H}^T\mathbf{H}\right)^{-1}\mathbf{H}^T\mathbf{x}$$

$$\hat{\mathbf{s}}_{LS} = \mathbf{H}\hat{\boldsymbol{\theta}}_{LS} = \mathbf{H}\left(\mathbf{H}^T\mathbf{H}\right)^{-1}\mathbf{H}^T\mathbf{x}$$

# Comparing LSE to MLE

**Model**                  **Estimate**

$$\mathbf{x} = \mathbf{H}\boldsymbol{\theta} + \mathbf{e}$$

No Probability Model Needed

$$\hat{\boldsymbol{\theta}}_{LS} = \left(\mathbf{H}^T\mathbf{H}\right)^{-1}\mathbf{H}^T\mathbf{x}$$

$$\mathbf{x} = \mathbf{H}\boldsymbol{\theta} + \mathbf{w}$$

PDF Gaussian, White

$$\hat{\theta}_{ML} = \left(\mathbf{H}^T\mathbf{H}\right)^{-1}\mathbf{H}^T\mathbf{x}$$

If you assume Gaussian & apply this… BUT you are WRONG … you at least get the LSE!

# Empirical Risk Minimization

We can generalize MLE by replacing the (conditional) log loss term, with any other loss function.

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{n=1}^{N} \ell(\boldsymbol{y}_n, \boldsymbol{\theta}; \boldsymbol{x}_n)$$

This is known as empirical risk minimization or ERM, since it is the expected loss where the expectation is taken wrt the empirical distribution.