

Regression: Probabilistic perspective

Machine Learning

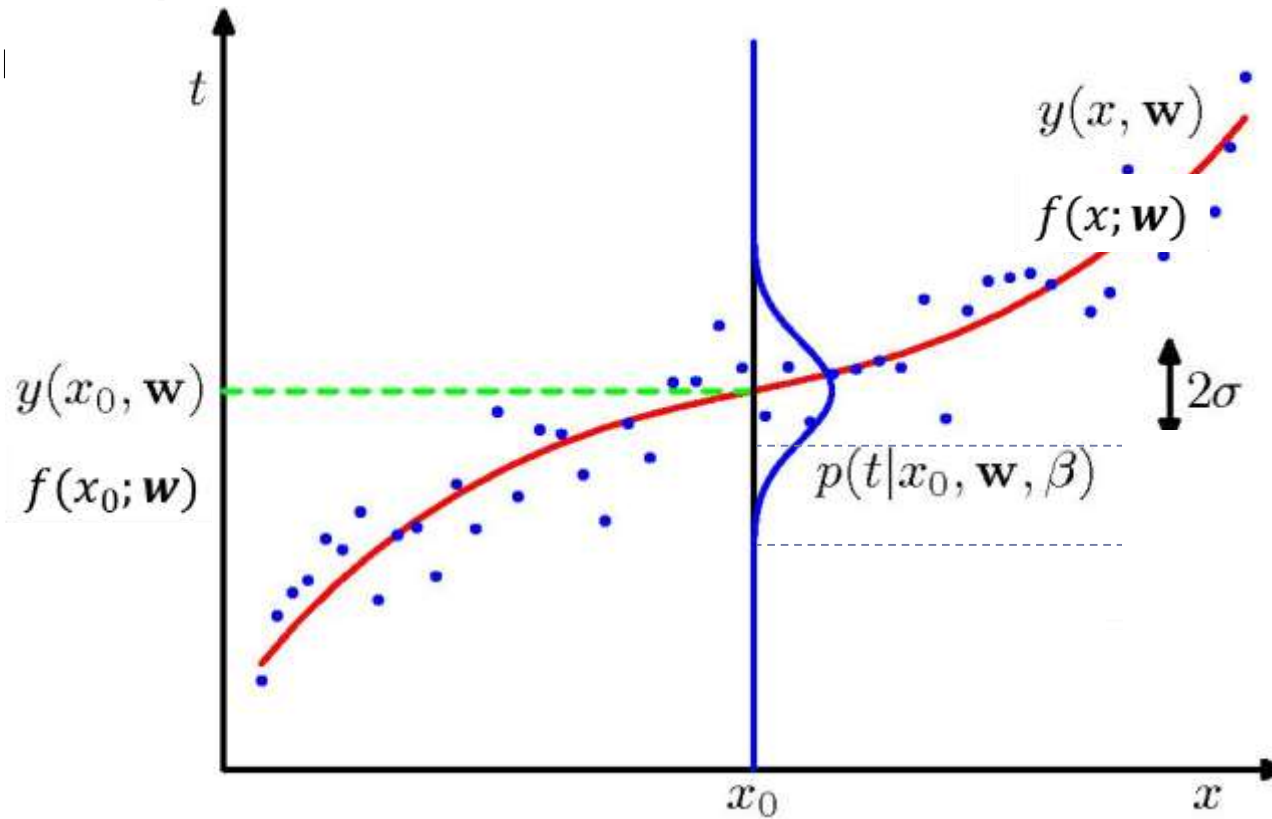
Hamid R Rabiee – Zahra Dehghanian
Spring 2025



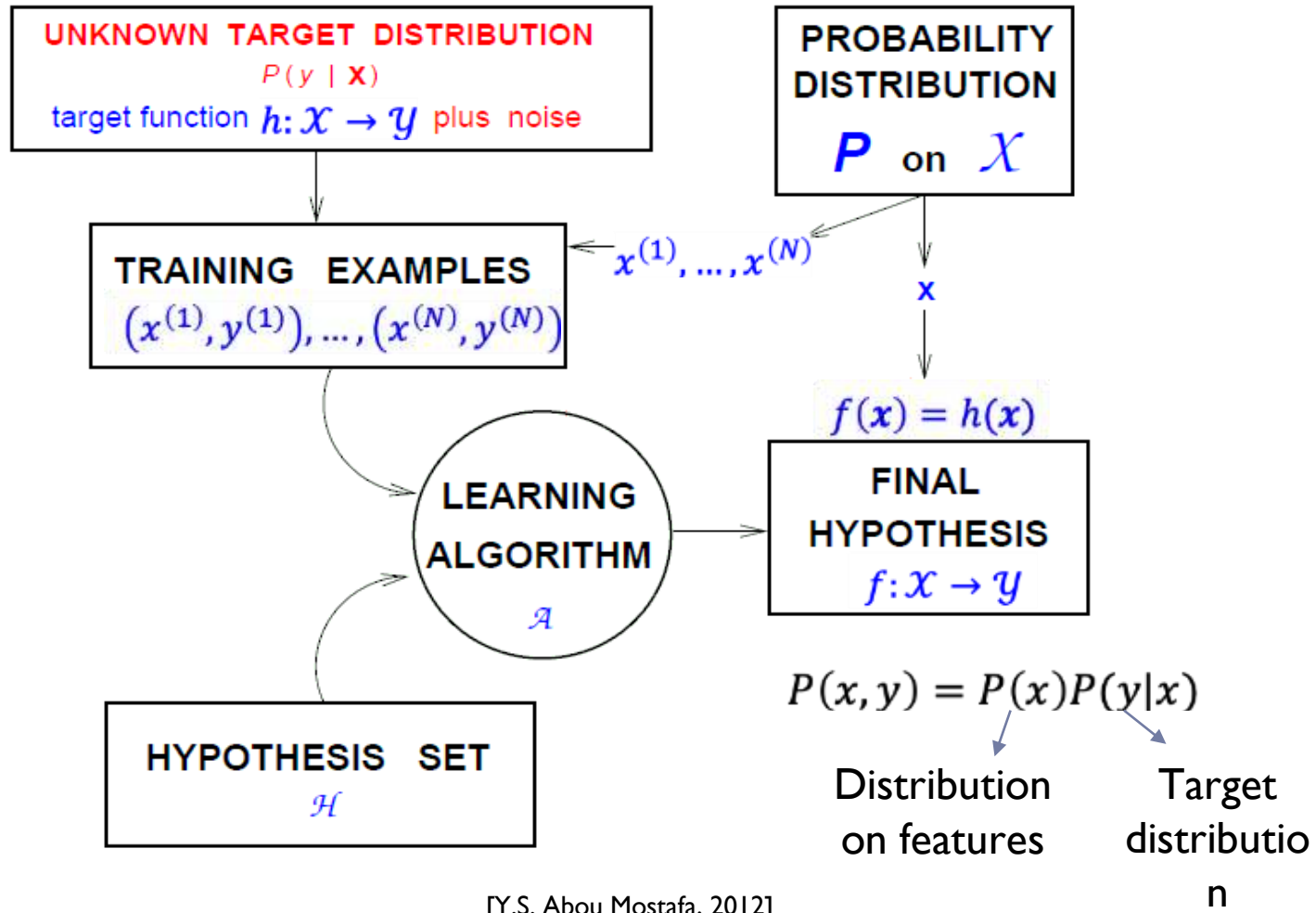
Sharif University
of Technology

Curve fitting: probabilistic perspective

- Describing uncertainty over value of target variable as a probability distribution
- Example



The learning diagram including noisy target



Curve fitting: probabilistic perspective (Example)

► Special case:

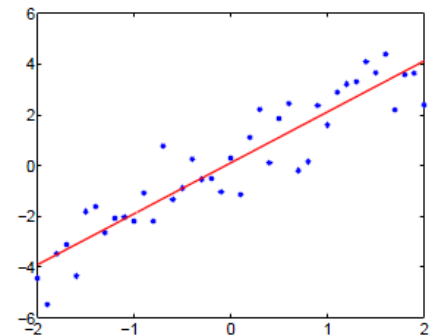
Observed output = function + noise

$$y = f(\mathbf{x}; \mathbf{w}) + \epsilon$$

$$\text{e.g., } \epsilon \sim N(0, \sigma^2)$$

$$y \sim N(f(\mathbf{x}; \mathbf{w}), \sigma^2)$$

- Noise: Whatever we cannot capture with our chosen family of functions



Curve fitting: probabilistic perspective (Example)

▶ Best regression

$$\mathbb{E}[y|\mathbf{x}] = E[f(\mathbf{x}; \mathbf{w}) + \epsilon] = f(\mathbf{x}; \mathbf{w})$$
$$\epsilon \sim N(0, \sigma^2)$$

- ▶ $f(\mathbf{x}; \mathbf{w})$ is trying to capture the mean of the observations y given the input \mathbf{x} :
- ▶ $\mathbb{E}[y|\mathbf{x}]$: conditional expectation of y given \mathbf{x}
 - ▶ evaluated according to the model (not according to the underlying distribution P)

Curve fitting using probabilistic estimation

- Maximum Likelihood (ML) estimation
- Maximum A Posteriori (MAP) estimation



Maximum likelihood estimation

- ▶ Given observations $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$
- ▶ Find the parameters that maximize the (conditional) likelihood of the outputs:

$$L(\mathcal{D}; \boldsymbol{\theta}) = p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = \prod_{i=1}^n p(y^{(i)}|\mathbf{x}^{(i)}, \boldsymbol{\theta})$$

$$\mathbf{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_d^{(1)} \\ 1 & x_1^{(2)} & \dots & x_d^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & \dots & x_d^{(n)} \end{bmatrix}$$

Maximum likelihood estimation (Cont'd)

□

$$y = f(\mathbf{x}; \mathbf{w}) + \epsilon, \quad \epsilon \sim N(0, \sigma^2)$$

- ▶ y given \mathbf{x} is normally distributed with mean $f(\mathbf{x}; \mathbf{w})$ and variance σ^2 :
 - ▶ we can also model the uncertainty in the predictions, not just the mean

$$p(y|\mathbf{x}, \mathbf{w}, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2\sigma^2} (y - f(\mathbf{x}; \mathbf{w}))^2\right\}$$



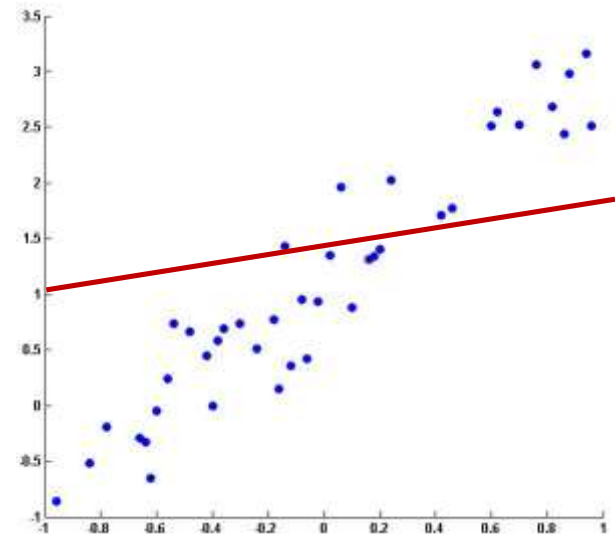
Maximum likelihood estimation (Cont'd)

▢ Example: univariate linear function

$$p(y|x, \mathbf{w}, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2\sigma^2} (y - w_0 - w_1x)^2\right\}$$

Why is this line a bad fit according to the likelihood criterion?

$p(y|x, \mathbf{w}, \sigma^2)$ for most of the points will be near zero (as they are far from this line)



Maximum likelihood estimation (Cont'd)

- ▶ Maximize the likelihood of the outputs (i.i.d):

$$L(\mathcal{D}; \mathbf{w}, \sigma^2) = \prod_{i=1}^N p(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}, \sigma^2)$$

$$\hat{\mathbf{w}} = \operatorname{argmax}_{\mathbf{w}} L(\mathcal{D}; \mathbf{w}, \sigma^2)$$

$$= \operatorname{argmax}_{\mathbf{w}} \prod_{i=1}^N p(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}, \sigma^2)$$



Maximum likelihood estimation (Cont'd)

- ▶ It is often easier (but equivalent) to try to maximize the log-likelihood:

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmax}} \ln p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \sigma^2)$$

$$\ln \prod_{i=1}^N p(y^{(i)}|\mathbf{x}^{(i)}, \mathbf{w}, \sigma^2) = \sum_{i=1}^N \ln \mathcal{N}(y^{(i)}|f(\mathbf{x}^{(i)}; \mathbf{w}), \sigma^2)$$

Maximum likelihood estimation (Cont'd)

- ▶ It is often easier (but equivalent) to try to maximize the log-likelihood:

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmax}} \ln p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \sigma^2)$$

$$\begin{aligned} \ln \prod_{i=1}^N p(y^{(i)}|\mathbf{x}^{(i)}, \mathbf{w}, \sigma^2) &= \sum_{i=1}^N \ln \mathcal{N}(y^{(i)}|f(\mathbf{x}^{(i)}; \mathbf{w}), \sigma^2) \\ &= -N \ln \sigma - \frac{N}{2} \ln 2\pi - \frac{1}{2\sigma^2} \underbrace{\sum_{i=1}^N (y^{(i)} - f(\mathbf{x}^{(i)}; \mathbf{w}))^2}_{\text{sum of squares error}} \end{aligned}$$

Maximum likelihood estimation (Cont'd)

- ▶ Maximizing log-likelihood (when we assume $y = f(\mathbf{x}; \mathbf{w}) + \epsilon$, $\epsilon \sim N(0, \sigma^2)$) is equivalent to minimizing SSE
- ▶ Let $\hat{\mathbf{w}}$ be the maximum likelihood (here least squares) setting of the parameters.



Maximum likelihood estimation (Cont'd)

- ▣ Maximizing log-likelihood (when we assume $y = f(\mathbf{x}; \mathbf{w}) + \epsilon$, $\epsilon \sim N(0, \sigma^2)$) is equivalent to minimizing SSE
- ▶ Let $\hat{\mathbf{w}}$ be the maximum likelihood (here least squares) setting of the parameters.
- ▶ What is the maximum likelihood estimate of σ^2 ?

$$\frac{\partial \log L(\mathcal{D}; \mathbf{w}, \sigma^2)}{\partial \sigma^2} = 0$$
$$\Rightarrow \hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - f(\mathbf{x}^{(i)}; \hat{\mathbf{w}}))^2$$

Mean squared prediction error

Maximum likelihood estimation (Cont'd)

- ▣ Generally, maximizing log-likelihood is equivalent to minimizing empirical loss when the loss is defined according to:

$$Loss\left(y^{(i)}, f(\mathbf{x}^{(i)}, \mathbf{w})\right) = -\ln p(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}, \boldsymbol{\theta})$$

- ▶ **Loss: negative log-probability**
 - ▶ More general distributions for $p(y|\mathbf{x})$ can be considered



Maximum A Posterior (MAP) estimation

▶ MAP:

- ▶ Given observations \mathcal{D}
- ▶ Find the parameters that maximize the probabilities of the parameters after observing the data (posterior probabilities):

$$\boldsymbol{\theta}_{MAP} = \max_{\boldsymbol{\theta}} p(\boldsymbol{\theta}|\mathcal{D})$$

Since $p(\boldsymbol{\theta}|\mathcal{D}) \propto p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})$

$$\boldsymbol{\theta}_{MAP} = \max_{\boldsymbol{\theta}} p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})$$

Maximum A Posterior (MAP) estimation

Given observations $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$

$$\max_{\mathbf{w}} p(\mathbf{w}|\mathbf{X}, \mathbf{y}) \propto p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})$$

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \alpha^2 \mathbf{I}) = \left(\frac{1}{\sqrt{2\pi}\alpha}\right)^{d+1} \exp\left\{-\frac{1}{2\alpha^2} \mathbf{w}^T \mathbf{w}\right\}$$

Maximum A Posterior (MAP) estimation

- ▶ Given observations $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$

$$\max_{\mathbf{w}} \ln p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \sigma^2)p(\mathbf{w})$$

$$\min_{\mathbf{w}} \frac{1}{\sigma^2} \sum_{i=1}^N (y^{(i)} - f(\mathbf{x}^{(i)}; \mathbf{w}))^2 + \frac{1}{\alpha^2} \mathbf{w}^T \mathbf{w}$$

- ▶ Equivalent to regularized SSE with $\lambda = \frac{\sigma^2}{\alpha^2}$

Feed back

- <https://forms.gle/vKRbyVVsWRKcZuqr8>



Resource

- C. Bishop, “Pattern Recognition and Machine Learning”, Chapter 3.3.
- Course CE-717, Dr. M.Soleymani

