

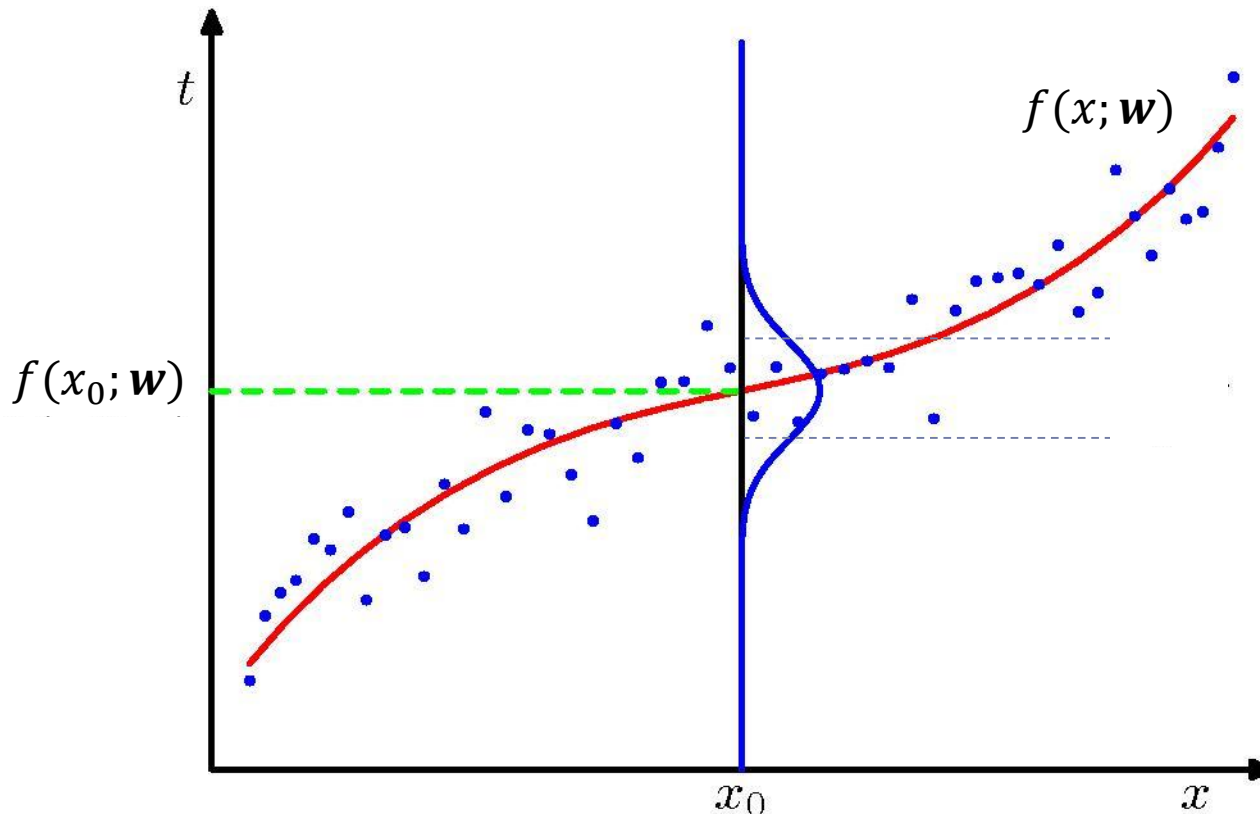
Regression and generalization

CE-717: Machine Learning
Sharif University of Technology

M. Soleymani
Fall 2016

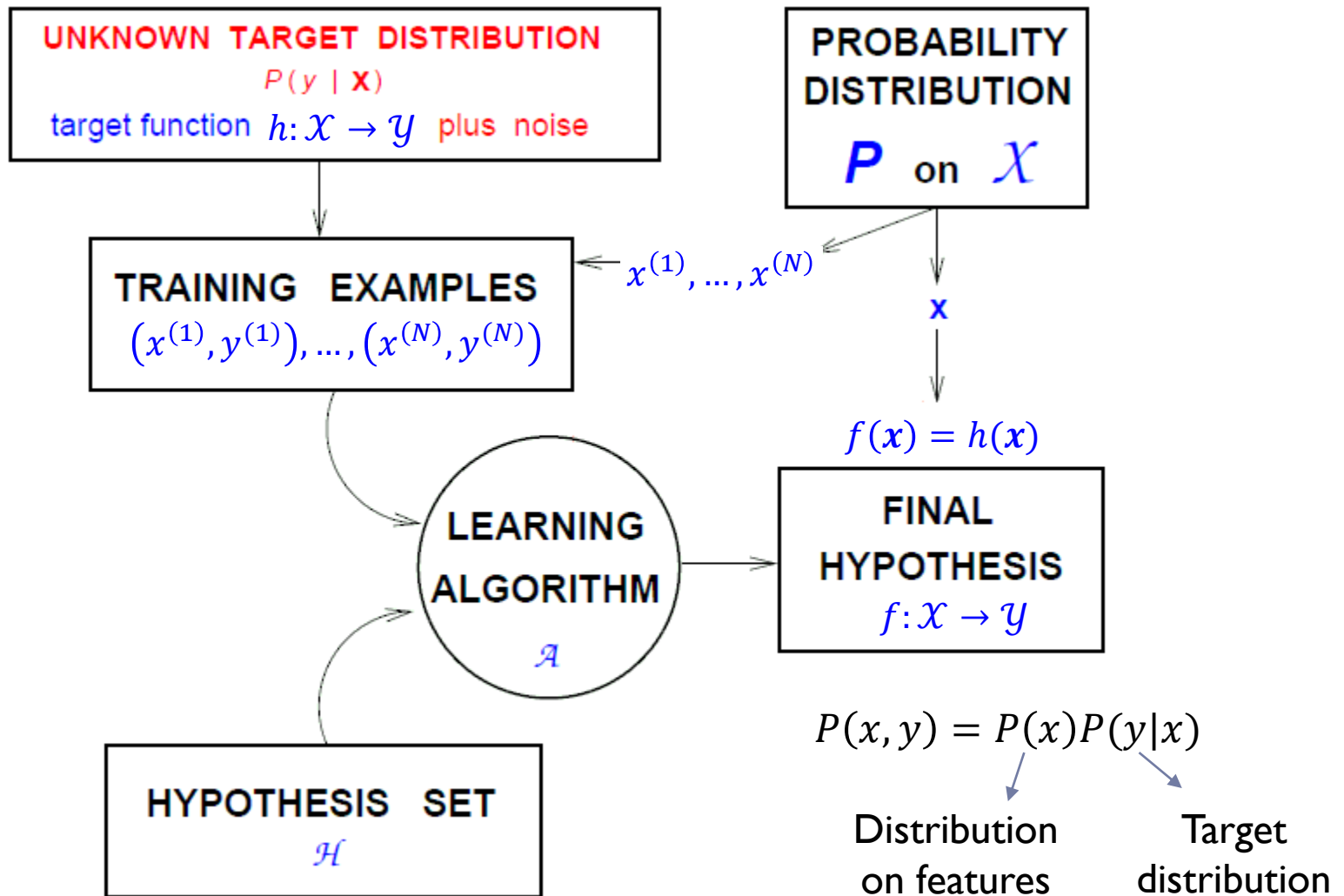
Curve fitting: probabilistic perspective

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



The learning diagram including noisy target

► Type



Curve fitting: probabilistic perspective (Example)

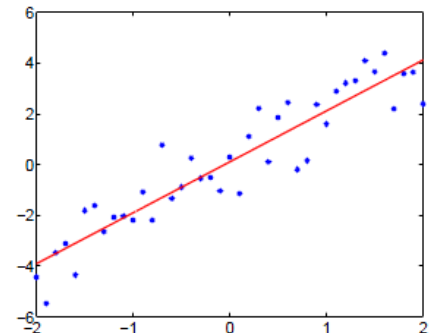
- ▶ Special case:

Observed output = function + noise

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

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

- ▶ 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
- ▶ Bayesian approach

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 model the uncertainty in the predictions, not just the mean

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

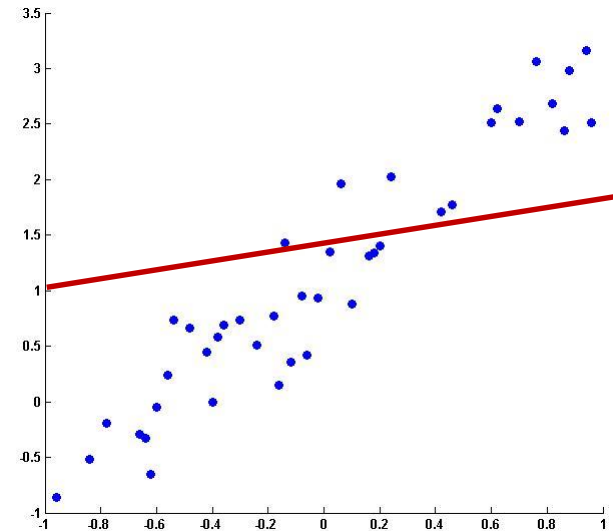
Maximum likelihood estimation (Cont'd)

- ▶ Example: univariate linear function

$$p(y|\mathbf{x}, \mathbf{w}, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2\sigma^2} (y - w_0 - w_1 x)^2\right\}$$

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

$p(y|\mathbf{x}, \mathbf{w}, \sigma^2)$ 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):

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

$$\begin{aligned} \hat{\mathbf{w}} &= \operatorname{argmax}_{\mathbf{w}} L(\mathcal{D}; \mathbf{w}, \sigma^2) \\ &= \operatorname{argmax}_{\mathbf{w}} \prod_{i=1}^n p(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}, \sigma^2) \end{aligned}$$

Maximum likelihood estimation (Cont'd)

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

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

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

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.
- ▶ 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 $\lambda = \frac{\sigma^2}{\alpha^2}$

Bayesian approach

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

$$p(y|\mathbf{x}, \mathcal{D}) = \int p(y|\mathbf{w}, \mathbf{x})p(\mathbf{w}|\mathcal{D})d\mathbf{w}$$

- ▶ Example of prior distribution

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \alpha^2 \mathbf{I})$$

Bayesian approach

- ▶ Given observations $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$
- ▶ Find the parameters that maximize the probabilities of observations

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

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \alpha^2 \mathbf{I})$$

$$p(\mathbf{w}|\mathcal{D}) \propto p(\mathcal{D}|\mathbf{w})p(\mathbf{w})$$

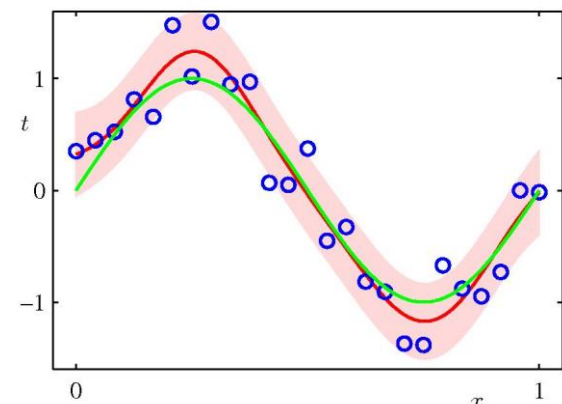
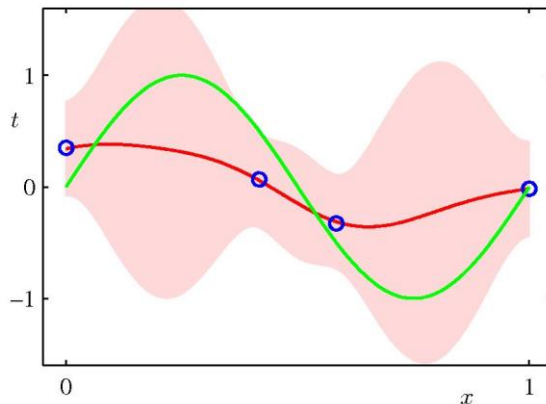
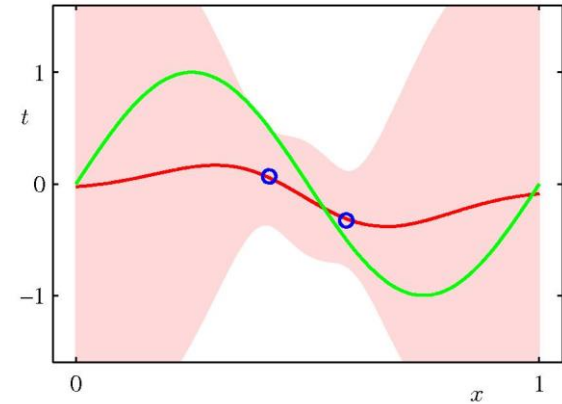
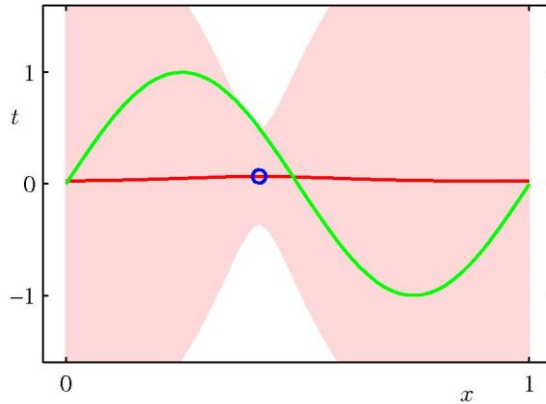
Predictive
distribution

$$p(y|\mathbf{x}, \mathcal{D}) = \int p(y|\mathbf{w}, \mathbf{x})p(\mathbf{w}|\mathcal{D})d\mathbf{w}$$

$$p(y|\mathbf{x}, \mathcal{D}) = N(\mathbf{m}_N^T \mathbf{x}, \sigma_N^2(\mathbf{x}))$$

Predictive distribution: example

- ▶ Example: Sinusoidal data, 9 Gaussian basis functions



Red curve shows the mean of the predictive distribution

Pink region spans one standard deviation either side of the mean

Predictive distribution: example

- ▶ Functions whose parameters are sampled from $p(\mathbf{w}|\mathcal{D})$

