
Probability Theory and Parameter Estimation II

Least Squares and Gauss



- **How** to solve the least square problem?
 - Carl Friedrich **Gauss** solution in 1794 (age 18)
 - **Why** solving the least square problem?
 - Carl Friedrich **Gauss** solution in 1822 (age 46)
 - Least square solution is optimal in the sense that it is the best linear unbiased estimator of the coefficients of the polynomials
 - **Assumptions:** errors have zero mean and equal variances
- http://en.wikipedia.org/wiki/Least_squares
-

Three Approaches

$p(\text{Parameters} \mid \text{Data})$

$p(\text{Data} \mid \text{Parameters})$

$p(\text{Parameters})$

$$p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}$$

posterior \propto likelihood \times prior

1. find parameters that maximize (log) likelihood
 2. find parameters that maximize posterior (MAP)
 3. find the posterior (fully Bayesian)
-

Maximum Likelihood I

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n|y(x_n, \mathbf{w}), \beta^{-1})$$

Maximize log likelihood

$$\ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = - \underbrace{\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2}_{\beta E(\mathbf{w})} + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$$

Surprise! Maximizing log likelihood is equivalent to minimizing sum of square error function!

Maximum Likelihood II

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | y(x_n, \mathbf{w}), \beta^{-1})$$

Maximize log likelihood

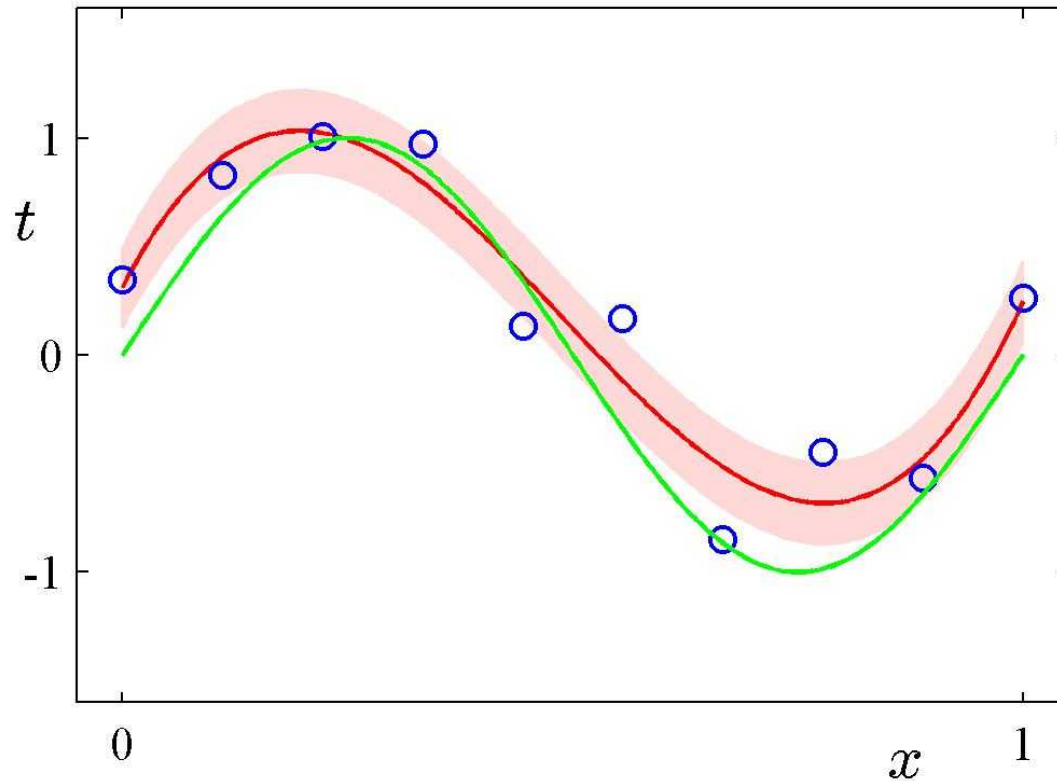
$$\ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = - \underbrace{\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2}_{\beta E(\mathbf{w})} + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$$

Determine \mathbf{w}_{ML} by minimizing sum-of-squares error, $E(\mathbf{w})$.

$$\frac{1}{\beta_{\text{ML}}} = \frac{1}{N} \sum_{n=1}^N \{y(x_n, \mathbf{w}_{\text{ML}}) - t_n\}^2$$

Predictive Distribution

$$p(t|x, \mathbf{w}_{\text{ML}}, \beta_{\text{ML}}) = \mathcal{N}(t|y(x, \mathbf{w}_{\text{ML}}), \beta_{\text{ML}}^{-1})$$



MAP: A Step towards Bayes

Maximum posterior

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

hyper-parameter

$$p(\mathbf{w}|\mathbf{x}, \mathbf{t}, \alpha, \beta) \propto p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta)p(\mathbf{w}|\alpha)$$

prior over parameters

likelihood

$$\beta\tilde{E}(\mathbf{w}) = \frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\alpha}{2}\mathbf{w}^T\mathbf{w}$$

Determine \mathbf{w}_{MAP} by minimizing regularized sum-of- $\tilde{E}(\mathbf{w})$ squares error.

MAP: A Step towards Bayes

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

$$p(\mathbf{w}|\mathbf{x}, \mathbf{t}, \alpha, \beta) \propto p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta)p(\mathbf{w}|\alpha)$$

$$\beta\tilde{E}(\mathbf{w}) = \frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\alpha}{2}\mathbf{w}^T\mathbf{w}$$

Determine \mathbf{w}_{MAP} by minimizing regularized sum-of- $\tilde{E}(\mathbf{w})$ squares error.

Surprise! Maximizing posterior is equivalent to minimizing regularized sum of square error function!

Three Approaches

$p(\text{Parameters} | \text{Data})$

$p(\text{Data} | \text{Parameters})$

$p(\text{Parameters})$

$$p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}$$

posterior \propto likelihood \times prior

1. find parameters that maximize (log) likelihood
2. find parameters that maximize posterior (MAP)
3. find the posterior (fully Bayesian)

$$p(t_0|X, x_0) = \int p(t_0|X, x_0, Y) p(Y|X, x_0) dY$$

Bayesian Curve Fitting

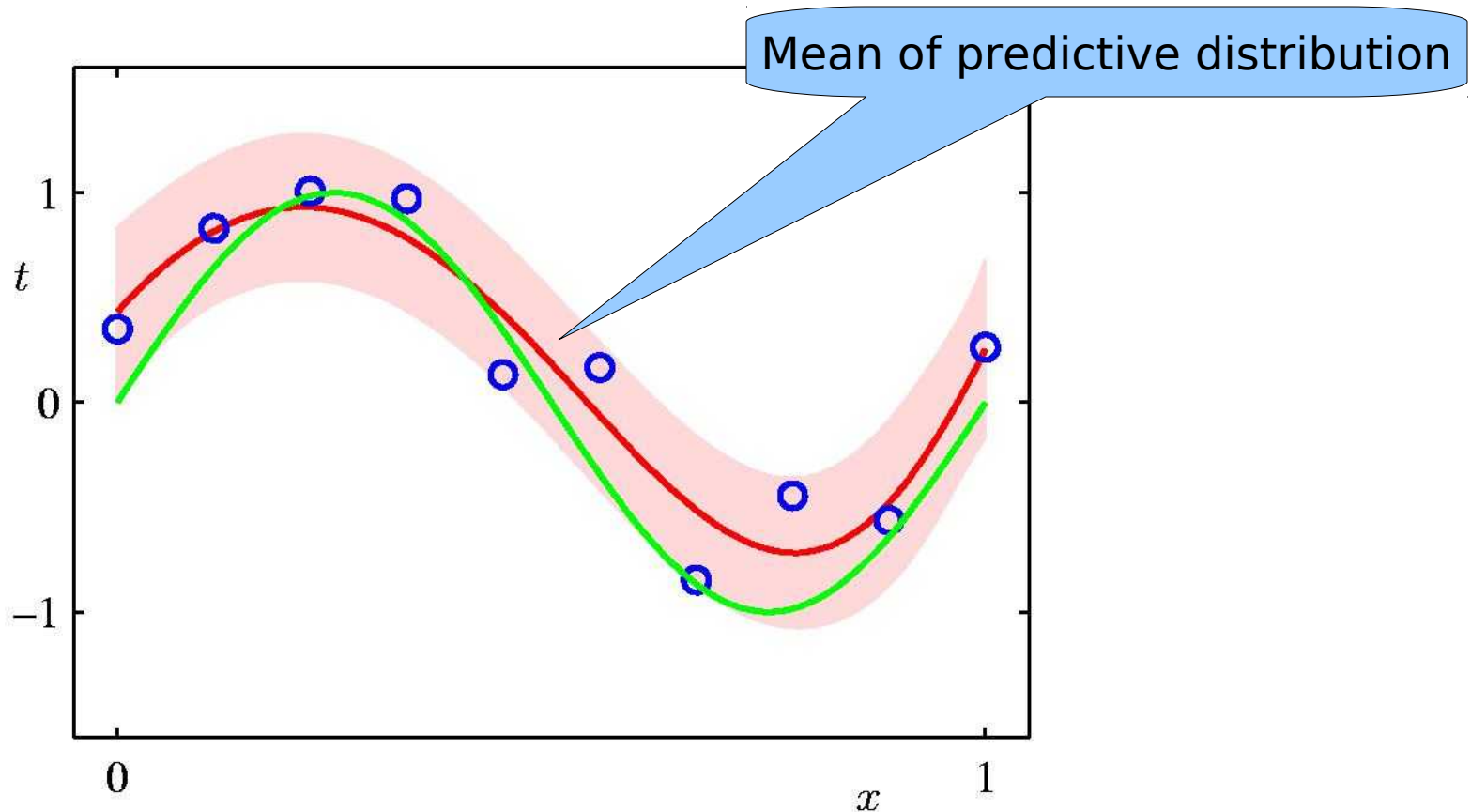
$$p(t|x, \mathbf{x}, \mathbf{t}) = \int p(t|x, \mathbf{w})p(\mathbf{w}|\mathbf{x}, \mathbf{t}) d\mathbf{w} = \mathcal{N}(t|m(x), s^2(x))$$

$$m(x) = \beta\phi(x)^T \mathbf{S} \sum_{n=1}^N \phi(x_n)t_n \quad s^2(x) = \beta^{-1} + \phi(x)^T \mathbf{S} \phi(x)$$

$$\mathbf{S}^{-1} = \alpha \mathbf{I} + \beta \sum_{n=1}^N \phi(x_n)\phi(x_n)^T \quad \phi(x_n) = (x_n^0, \dots, x_n^M)^T$$

Bayesian Predictive Distribution

$$p(t|x, \mathbf{x}, \mathbf{t}) = \mathcal{N}(t|m(x), s^2(x))$$



Review

$$\text{posterior} \propto \text{likelihood} \times \text{prior}$$

1. find parameters that maximize (log) likelihood	can lead to over-fitting (yields parameters)
2. find parameters that maximize posterior (MAP)	avoids over-fitting (yields parameters)
3. find the posterior (fully Bayesian)	yields distribution

Review

posterior \propto likelihood \times prior

1. find parameters that maximize (log) likelihood

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2$$

2. find parameters that maximize posterior (MAP)

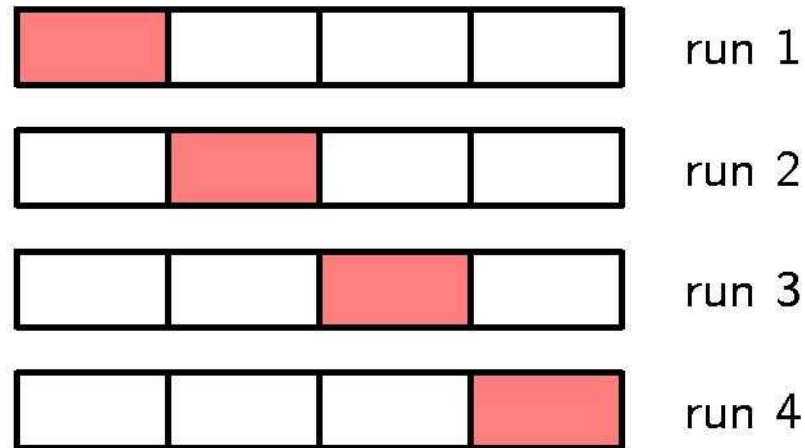
$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

3. find the posterior (fully Bayesian)

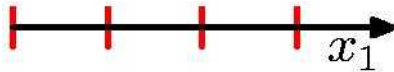
yields distribution

Model Selection

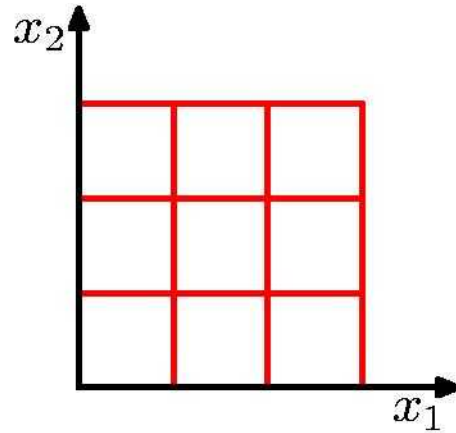
Cross-Validation



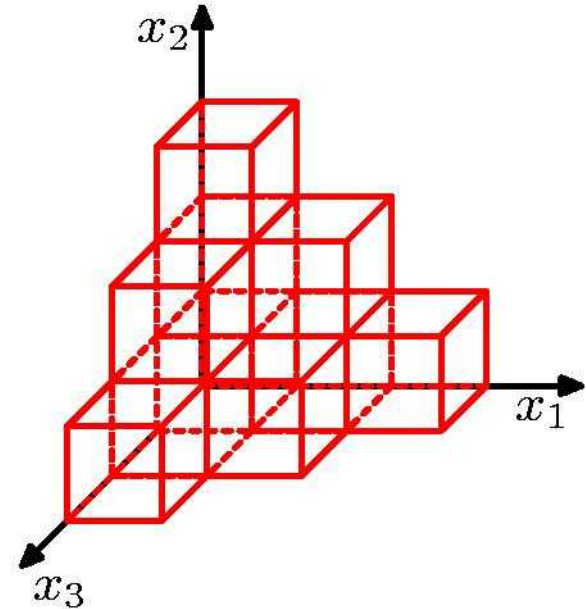
Curse of Dimensionality



$D = 1$

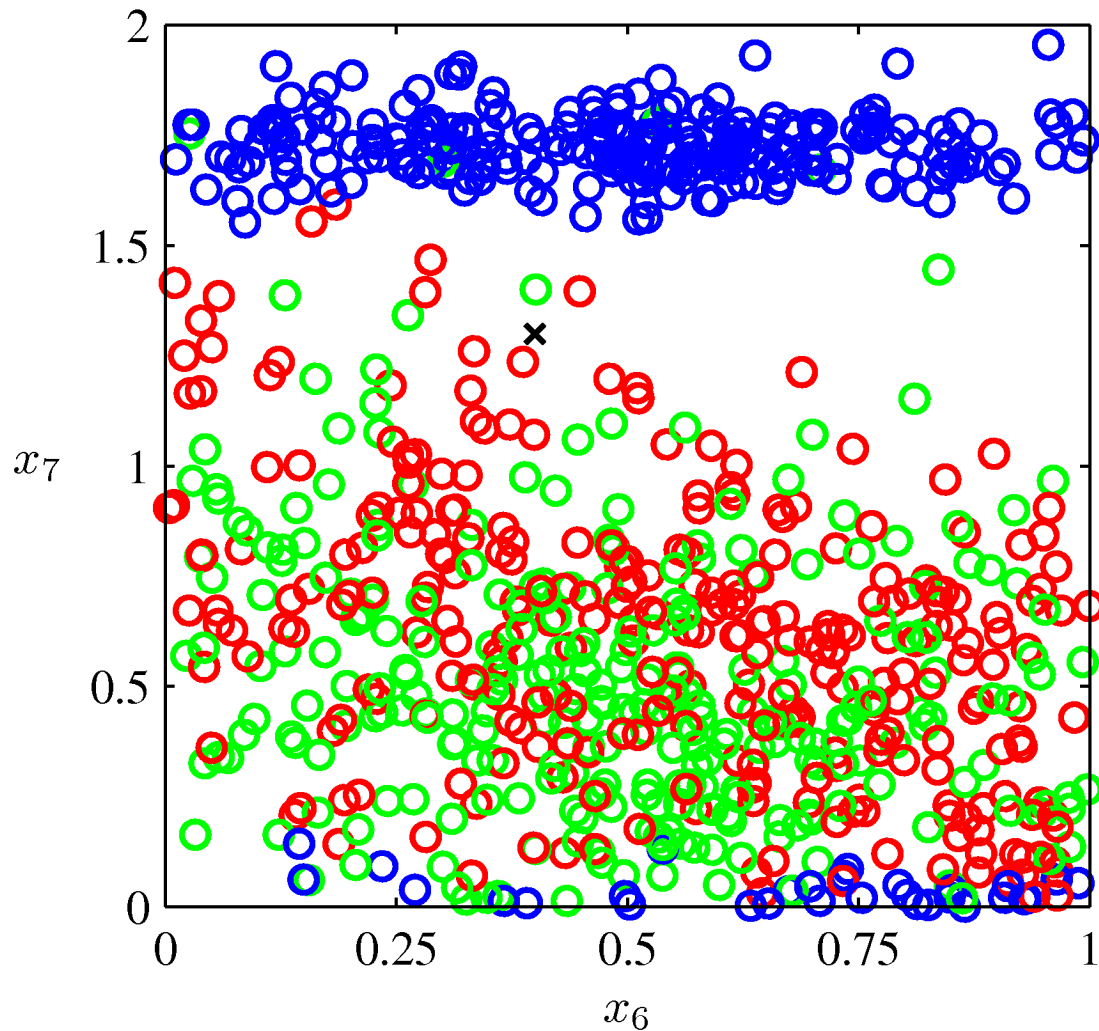


$D = 2$

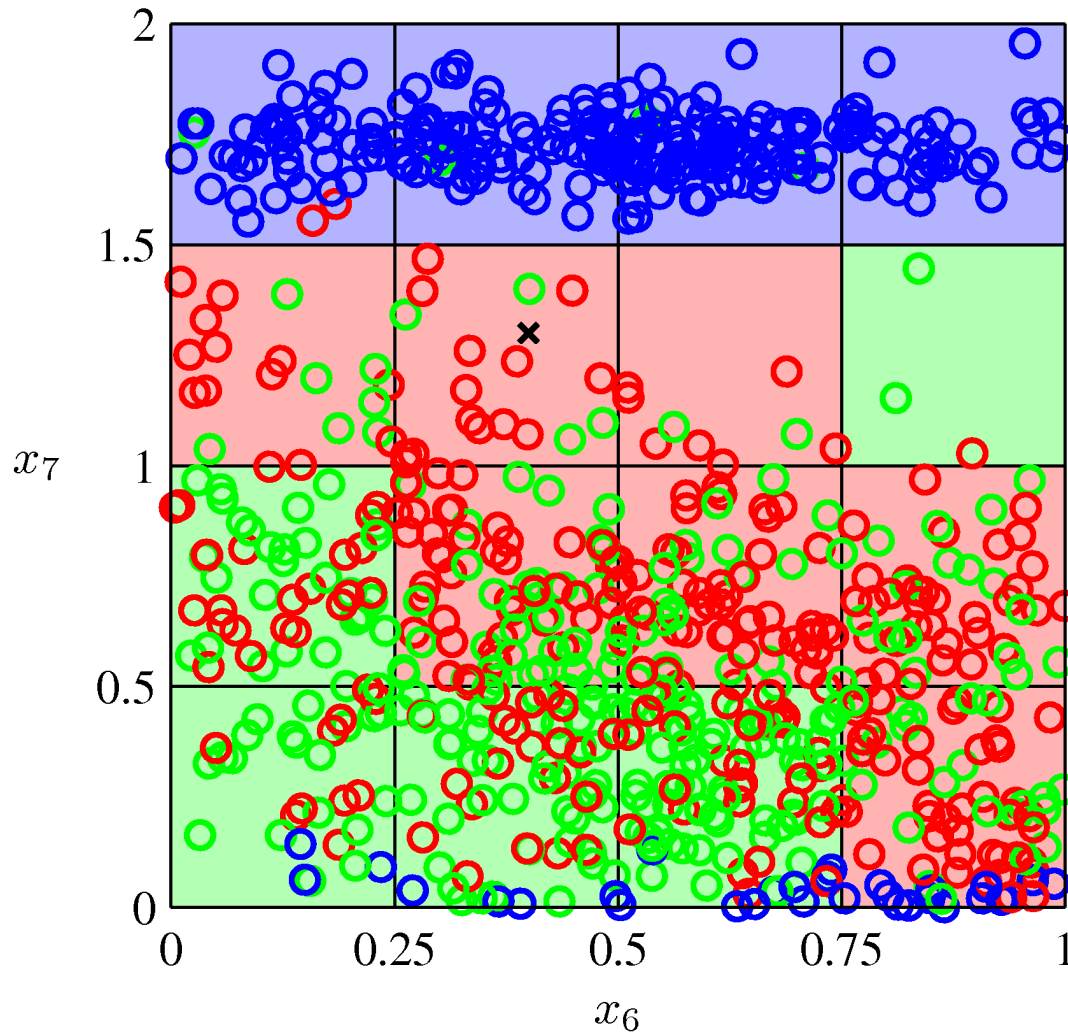


$D = 3$

Curse of Dimensionality



Curse of Dimensionality



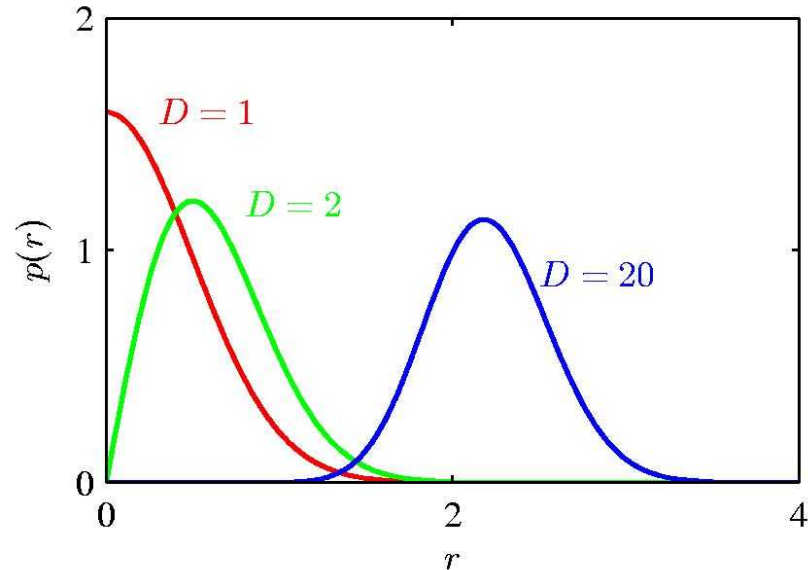
Curse of Dimensionality

Polynomial curve fitting, M

$= 3$

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{i=1}^D w_i x_i + \sum_{i=1}^D \sum_{j=1}^D w_{ij} x_i x_j + \sum_{i=1}^D \sum_{j=1}^D \sum_{k=1}^D w_{ijk} x_i x_j x_k$$

Gaussian Densities in
higher dimensions



Decision Theory

Inference step

Determine either $p(t|\mathbf{x})$ or $p(\mathbf{x}, t)$.



regression

Decision step

For given \mathbf{x} , determine optimal a (action).

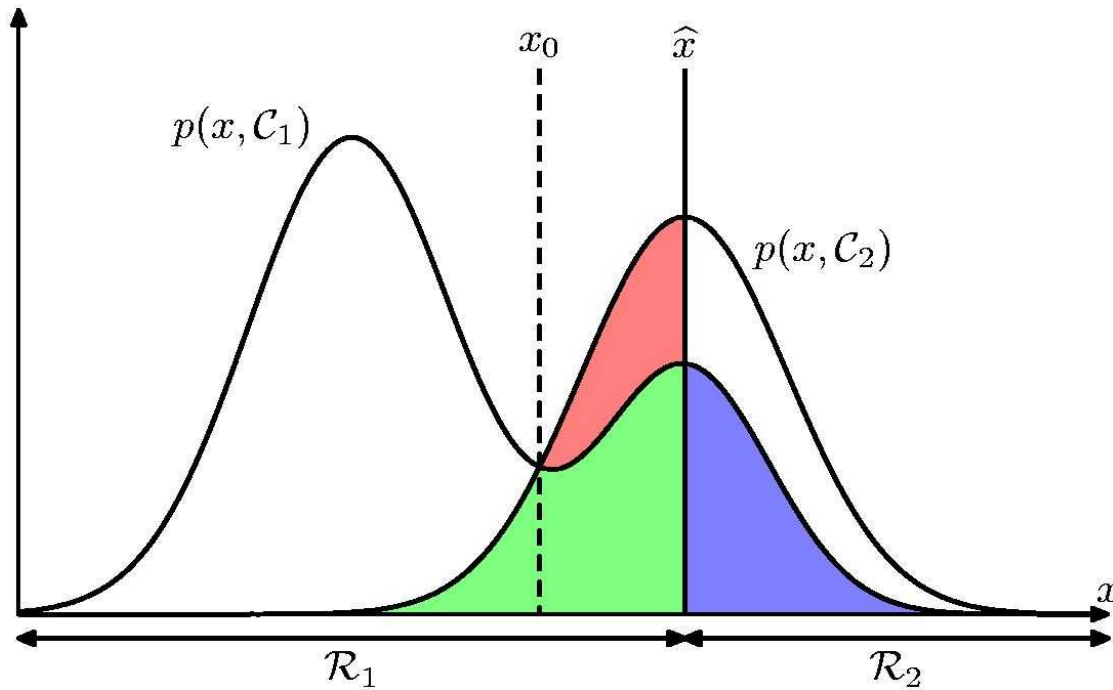


classification

$$p(\text{cancer}|\text{image}) = \frac{p(\text{image}|\text{cancer}) p(\text{cancer})}{p(\text{image})}$$

To minimize misclassification: maximize posterior

Minimum Misclassification Rate



$$\begin{aligned} p(\text{mistake}) &= p(\mathbf{x} \in \mathcal{R}_1, \mathcal{C}_2) + p(\mathbf{x} \in \mathcal{R}_2, \mathcal{C}_1) \\ &= \int_{\mathcal{R}_1} p(\mathbf{x}, \mathcal{C}_2) d\mathbf{x} + \int_{\mathcal{R}_2} p(\mathbf{x}, \mathcal{C}_1) d\mathbf{x}. \end{aligned}$$

Minimum Expected Loss

Example: classify medical images as 'cancer' or 'normal'

		Decision	
		cancer	normal
Loss function	Truth cancer	0	1000
	Truth normal	1	0

Minimum Expected Loss

elements in region j

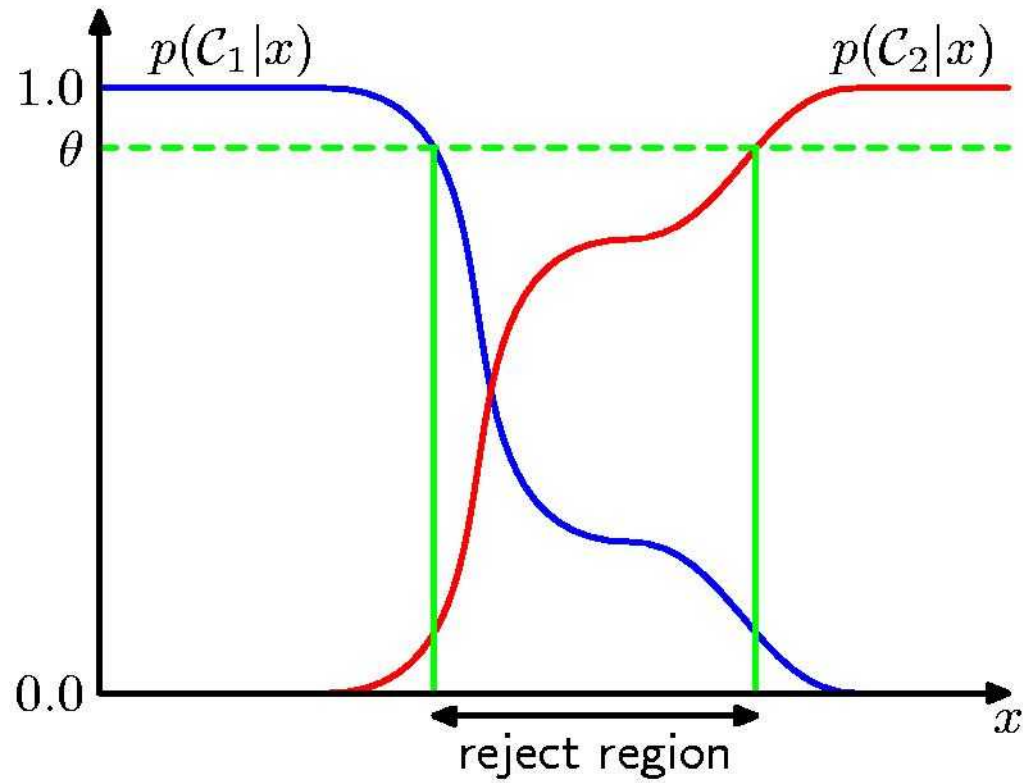
$$\mathbb{E}[L] = \sum_k \sum_j \int_{\mathcal{R}_j} L_{kj} p(\mathbf{x}, \mathcal{C}_k) d\mathbf{x}$$

real class is k

Regions \mathcal{R}_j are chosen to minimize

$$\sum_k L_{kj} p(\mathcal{C}_k | \mathbf{x})$$

Reject Option



Why Separate Inference and Decision?

- Minimizing risk (loss matrix may change over time)
 - Reject option
 - Unbalanced class priors
 - Combining models
-

Decision Theory for Regression

Inference step

Determine $p(\mathbf{x}, t)$.

Decision step

For given \mathbf{x} , make optimal prediction, $y(\mathbf{x})$, for t .

Loss function: $\mathbb{E}[L] = \iint L(t, y(\mathbf{x}))p(\mathbf{x}, t) d\mathbf{x} dt$

The Squared Loss Function

$$\mathbb{E}[L] = \iint \{y(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) \, d\mathbf{x} \, dt$$

$$\begin{aligned} \{y(\mathbf{x}) - t\}^2 &= \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}] + \mathbb{E}[t|\mathbf{x}] - t\}^2 \\ &= \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}^2 + 2\{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}\{\mathbb{E}[t|\mathbf{x}] - t\} + \{\mathbb{E}[t|\mathbf{x}] - t\}^2 \end{aligned}$$

$$\mathbb{E}[L] = \int \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}^2 p(\mathbf{x}) \, d\mathbf{x} + \int \text{var} [t|\mathbf{x}] p(\mathbf{x}) \, d\mathbf{x}$$

$$y(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}]$$



As expected