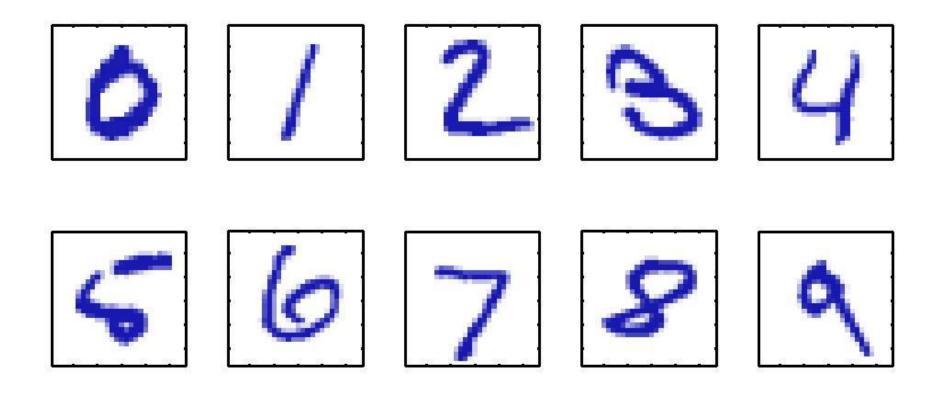
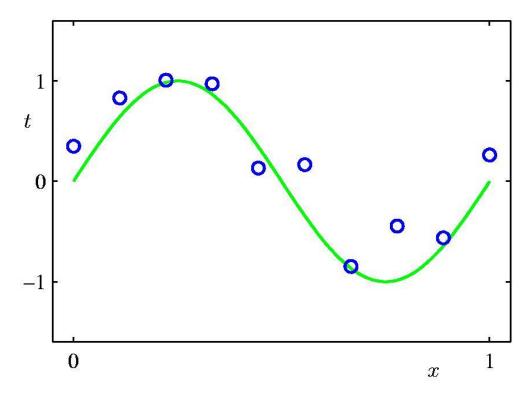


Example

Handwritten Digit Recognition

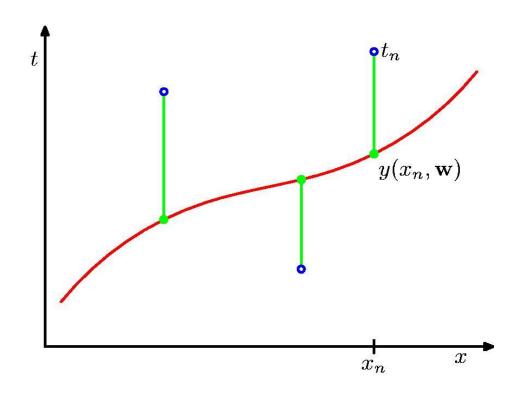


Polynomial Curve Fitting



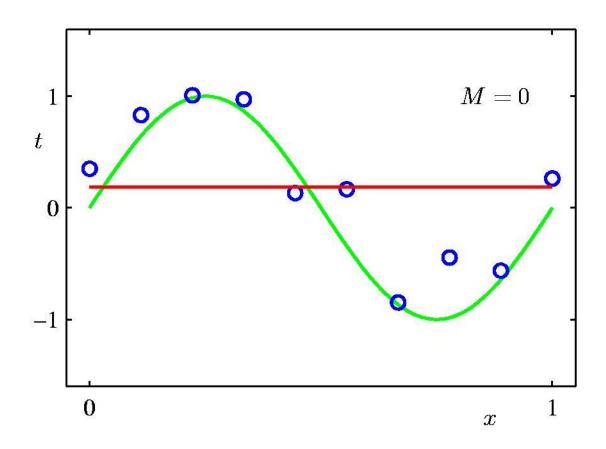
$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \ldots + w_M x^M = \sum_{j=0}^{M} w_j x^j$$

Sum-of-Squares Error Function

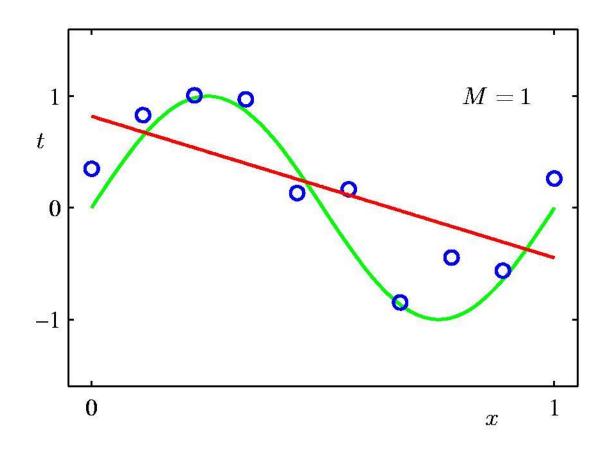


$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \{y(x_n, \mathbf{w}) - t_n\}^2$$

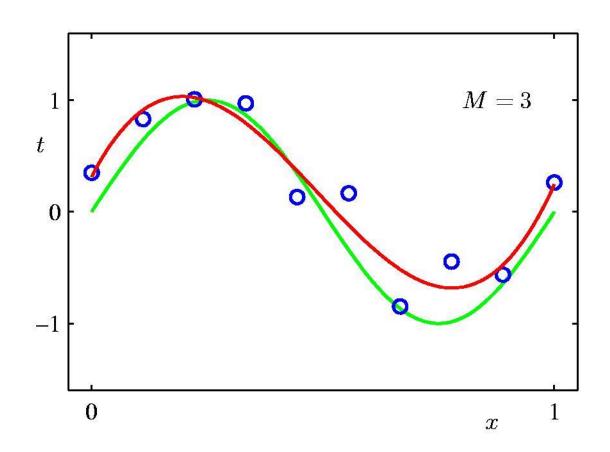
Oth Order Polynomial



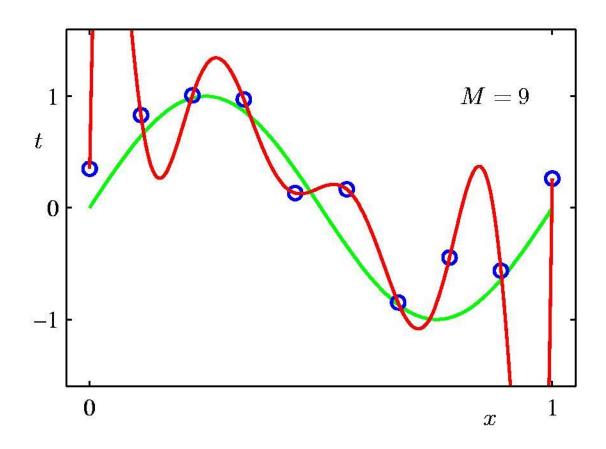
1st Order Polynomial



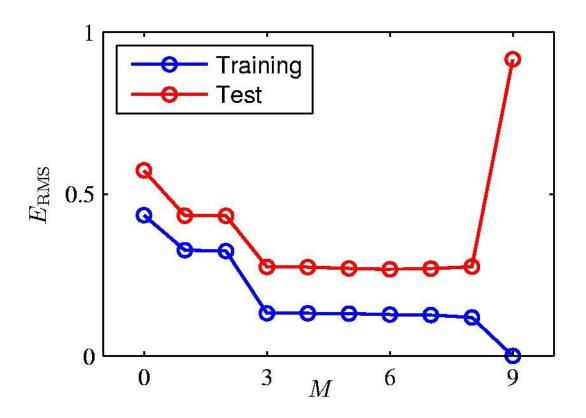
3rd Order Polynomial



9th Order Polynomial



Over-fitting



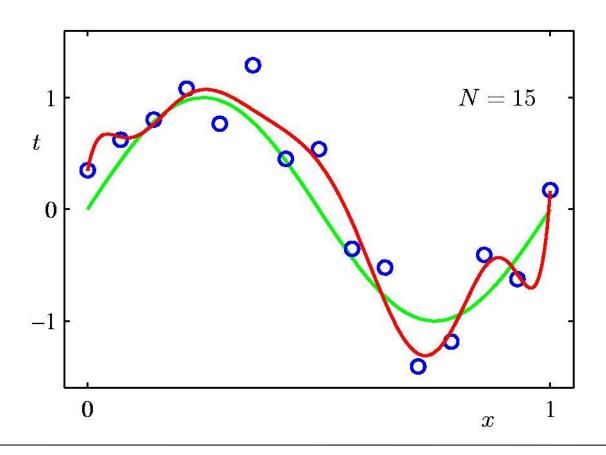
Root-Mean-Square (RMS) Error: $E_{\rm RMS} = \sqrt{2E(\mathbf{w}^\star)/N}$

Polynomial Coefficients

	M=0	M = 1	M = 3	M = 9
$\overline{w_0^{\star}}$	0.19	0.82	0.31	0.35
w_1^{\star}		-1.27	7.99	232.37
w_2^{\star}			-25.43	-5321.83
w_3^{\star}			17.37	48568.31
w_4^{\star}				-231639.30
w_5^{\star}				640042.26
w_6^{\star}				-1061800.52
w_7^\star				1042400.18
w_8^{\star}				-557682.99
w_9^{\star}				125201.43

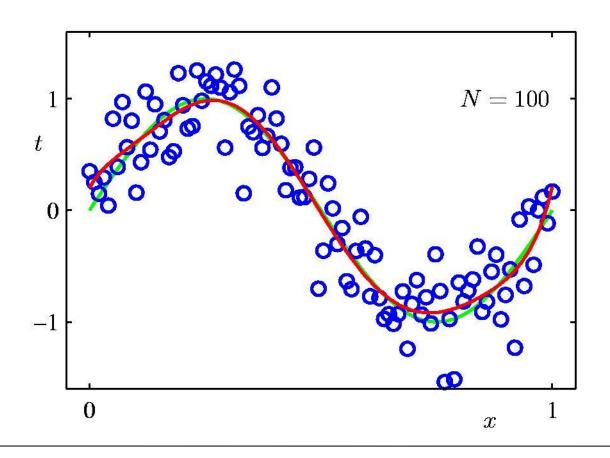
Data Set Size: N=15

9th Order Polynomial



Data Set Size: N = 100

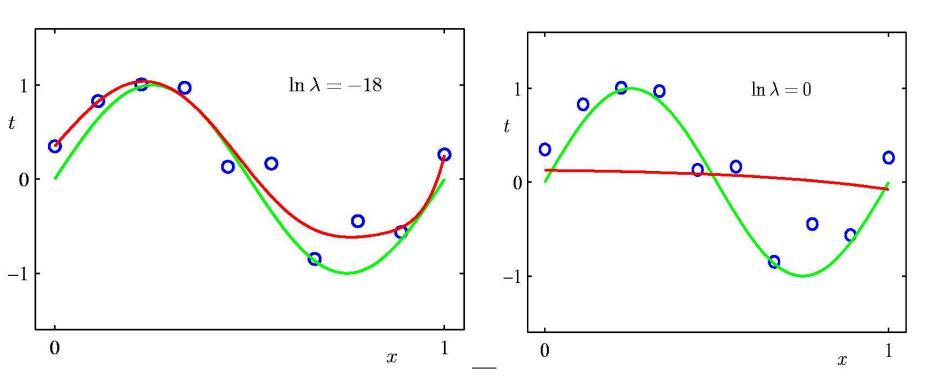
9th Order Polynomial



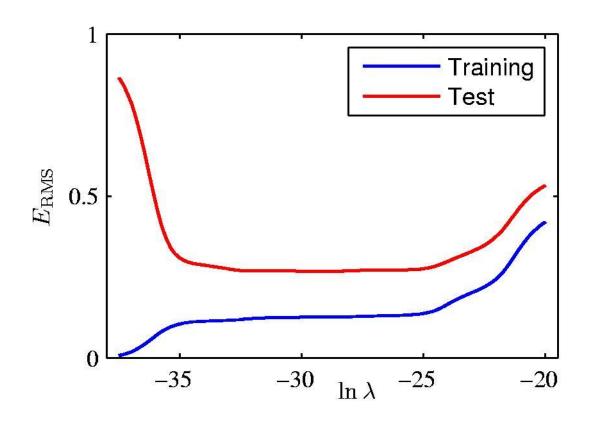
Regularization

Penalize large coefficient values

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



Regularization: $E_{\rm RMS}$ vs. $\ln \lambda$

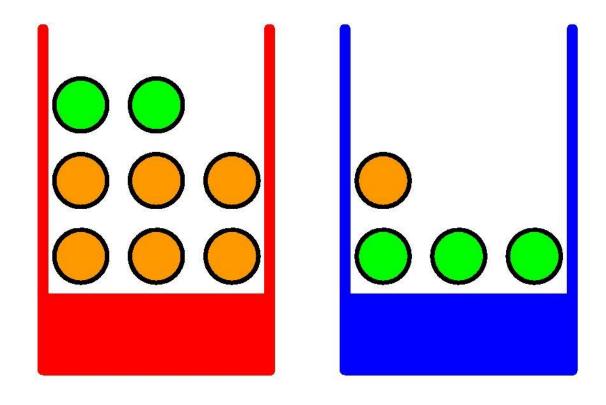


Polynomial Coefficients

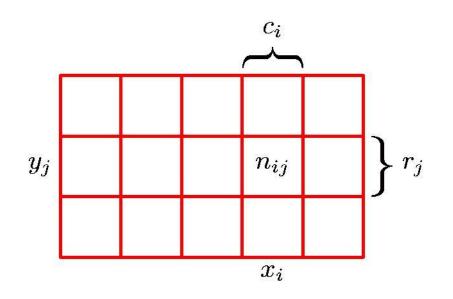
	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
w_0^{\star}	0.35	0.35	0.13
w_1^{\star}	232.37	4.74	-0.05
w_2^{\star}	-5321.83	-0.77	-0.06
w_3^{\star}	48568.31	-31.97	-0.05
w_4^{\star}	-231639.30	-3.89	-0.03
w_5^{\star}	640042.26	55.28	-0.02
w_6^{\star}	-1061800.52	41.32	-0.01
w_7^{\star}	1042400.18	-45.95	-0.00
w_8^\star	-557682.99	-91.53	0.00
w_9^{\star}	125201.43	72.68	0.01

Probability Theory

Apples and Oranges



Probability Theory



Marginal Probability

$$p(X = x_i) = \frac{c_i}{N}.$$

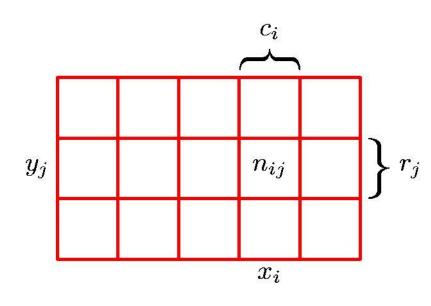
Joint Probability

$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N}$$

Conditional Probability

$$p(Y = y_j | X = x_i) = \frac{n_{ij}}{c_i}$$

Probability Theory



Sum Rule

$$r_j$$
 $p(X = x_i) = \frac{c_i}{N} = \frac{1}{N} \sum_{j=1}^{L} n_{ij}$
= $\sum_{j=1}^{L} p(X = x_i, Y = y_j)$

Product Rule

$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N} = \frac{n_{ij}}{c_i} \cdot \frac{c_i}{N}$$
$$= p(Y = y_j | X = x_i) p(X = x_i)$$

The Rules of Probability

Sum Rule

$$p(X) = \sum_{Y} p(X, Y)$$

Product Rule

$$p(X,Y) = p(Y|X)p(X)$$

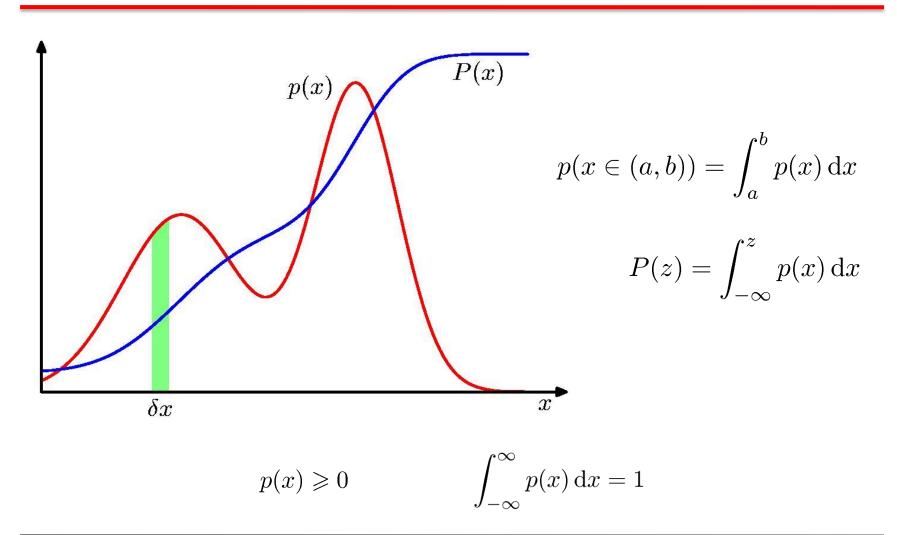
Bayes' Theorem

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

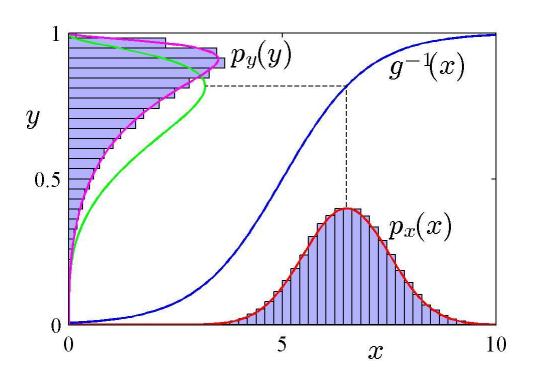
$$p(X) = \sum_{Y} p(X|Y)p(Y)$$

posterior ∞ likelihood × prior

Probability Densities



Transformed Densities



$$p_y(y) = p_x(x) \left| \frac{\mathrm{d}x}{\mathrm{d}y} \right|$$

= $p_x(g(y)) |g'(y)|$

Expectations

$$\mathbb{E}[f] = \sum_{x} p(x)f(x)$$

$$\mathbb{E}[f] = \int p(x)f(x) \, \mathrm{d}x$$

$$\mathbb{E}_x[f|y] = \sum_x p(x|y)f(x)$$

Conditional Expectation (discrete)

$$\mathbb{E}[f] \simeq \frac{1}{N} \sum_{n=1}^{N} f(x_n)$$

Approximate Expectation (discrete and continuous)

Variances and Covariances

$$\operatorname{var}[f] = \mathbb{E}\left[\left(f(x) - \mathbb{E}[f(x)]\right)^{2}\right] = \mathbb{E}[f(x)^{2}] - \mathbb{E}[f(x)]^{2}$$

$$cov[x, y] = \mathbb{E}_{x,y} [\{x - \mathbb{E}[x]\} \{y - \mathbb{E}[y]\}]$$

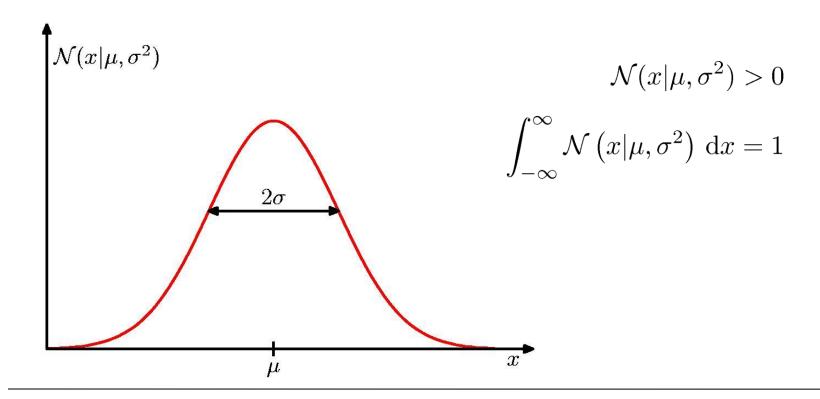
$$= \mathbb{E}_{x,y} [xy] - \mathbb{E}[x]\mathbb{E}[y]$$

$$cov[\mathbf{x}, \mathbf{y}] = \mathbb{E}_{\mathbf{x},\mathbf{y}} [\{\mathbf{x} - \mathbb{E}[\mathbf{x}]\} \{\mathbf{y}^{\mathrm{T}} - \mathbb{E}[\mathbf{y}^{\mathrm{T}}]\}]$$

$$= \mathbb{E}_{\mathbf{x},\mathbf{y}} [\mathbf{x}\mathbf{y}^{\mathrm{T}}] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{y}^{\mathrm{T}}]$$

The Gaussian Distribution

$$\mathcal{N}(x|\mu,\sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\}$$



Gaussian Mean and Variance

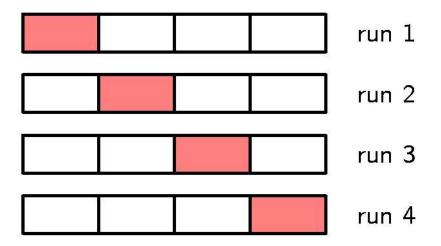
$$\mathbb{E}[x] = \int_{-\infty}^{\infty} \mathcal{N}(x|\mu, \sigma^2) x \, \mathrm{d}x = \mu$$

$$\mathbb{E}[x^2] = \int_{-\infty}^{\infty} \mathcal{N}(x|\mu, \sigma^2) x^2 dx = \mu^2 + \sigma^2$$

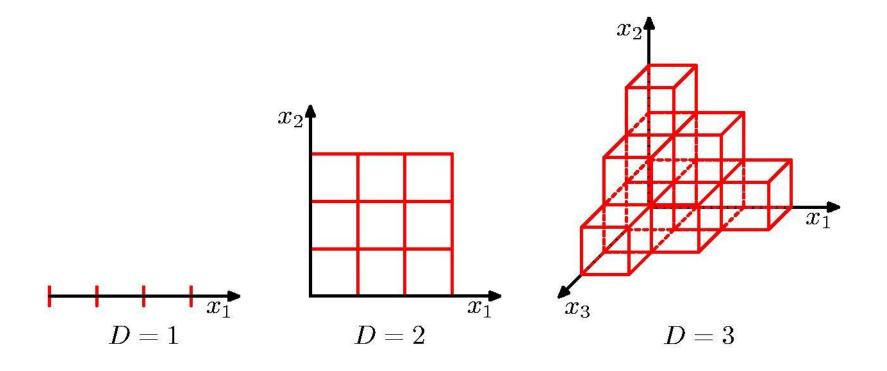
$$var[x] = \mathbb{E}[x^2] - \mathbb{E}[x]^2 = \sigma^2$$

Model Selection

Cross-Validation



Curse of Dimensionality



Decision Theory

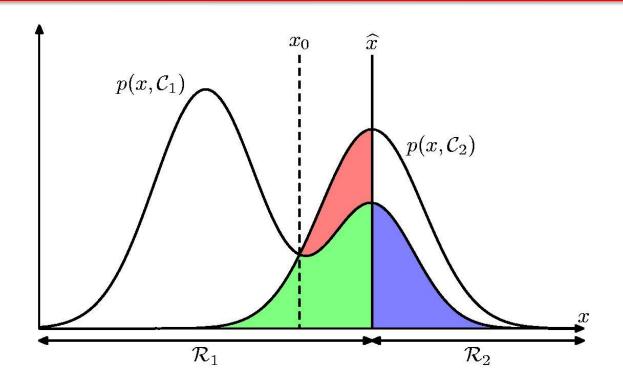
Inference step

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

Decision step

For given x, determine optimal t.

Minimum Misclassification Rate



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

Minimum Expected Loss

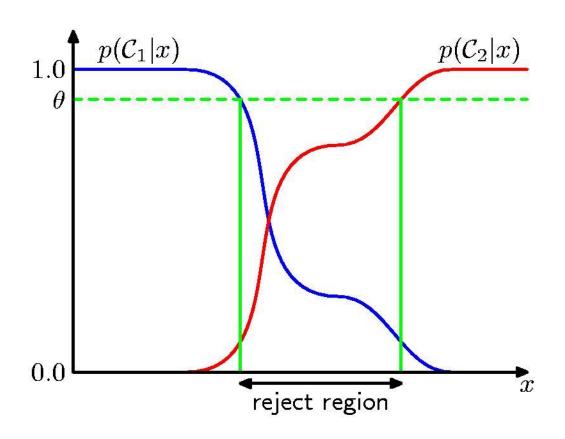
Example: classify medical images as 'cancer' or 'normal'

$$\mathbb{E}[L] = \sum_{k} \sum_{j} \int_{\mathcal{R}_{j}} L_{kj} p(\mathbf{x}, \mathcal{C}_{k}) d\mathbf{x}$$

Regions \mathcal{R}_i are chosen to minimize

$$\mathbb{E}[L] = \sum_{k} L_{kj} p(\mathcal{C}_k | \mathbf{x})$$

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

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

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

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

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

Generative vs Discriminative

Generative approach:

Model
$$p(t, \mathbf{x}) = p(\mathbf{x}|t)p(t)$$

Use Bayes' theorem $p(t|\mathbf{x}) = \frac{p(\mathbf{x}|t)p(t)}{p(\mathbf{x})}$

Discriminative approach:

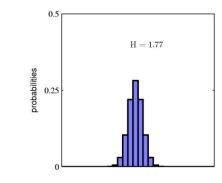
Model $p(t|\mathbf{x})$ directly

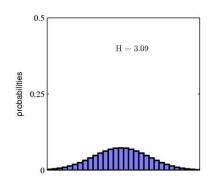
Entropy

$$H[x] = -\sum_{x} p(x) \log_2 p(x)$$

Important quantity in

- coding theory
- statistical physics
- machine learning





Conditional Entropy

$$H[\mathbf{y}|\mathbf{x}] = -\iint p(\mathbf{y}, \mathbf{x}) \ln p(\mathbf{y}|\mathbf{x}) \, d\mathbf{y} \, d\mathbf{x}$$

$$H[\mathbf{x}, \mathbf{y}] = H[\mathbf{y}|\mathbf{x}] + H[\mathbf{x}]$$

The Kullback-Leibler Divergence

$$KL(p||q) = -\int p(\mathbf{x}) \ln q(\mathbf{x}) d\mathbf{x} - \left(-\int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x}\right)$$
$$= -\int p(\mathbf{x}) \ln \left\{\frac{q(\mathbf{x})}{p(\mathbf{x})}\right\} d\mathbf{x}$$

$$\mathrm{KL}(p\|q) \simeq \frac{1}{N} \sum_{n=1}^{N} \left\{ -\ln q(\mathbf{x}_n|\boldsymbol{\theta}) + \ln p(\mathbf{x}_n) \right\}$$

$$KL(p||q) \geqslant 0$$
 $KL(p||q) \not\equiv KL(q||p)$

Mutual Information

$$I[\mathbf{x}, \mathbf{y}] \equiv KL(p(\mathbf{x}, \mathbf{y}) || p(\mathbf{x}) p(\mathbf{y}))$$

$$= -\iint p(\mathbf{x}, \mathbf{y}) \ln \left(\frac{p(\mathbf{x}) p(\mathbf{y})}{p(\mathbf{x}, \mathbf{y})} \right) d\mathbf{x} d\mathbf{y}$$

$$I[\mathbf{x}, \mathbf{y}] = H[\mathbf{x}] - H[\mathbf{x}|\mathbf{y}] = H[\mathbf{y}] - H[\mathbf{y}|\mathbf{x}]$$