

# Linear Regression and Classification

Wenting Tu

SHUFE, SIME

Machine Learning and Deep Learning

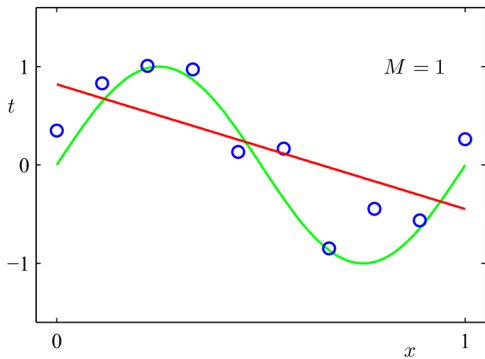
Course No. 1638

# Outline

Linear Regression

Linear Classification

# Illustration



$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + w_3x^3$$

## Linear basis function models

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x})$$

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) \quad \phi_0(\mathbf{x}) = 1$$

$$\phi_j(x) = x^j$$

$$\phi_j(x) = \exp \left\{ -\frac{(x - \mu_j)^2}{2s^2} \right\}$$

$$\phi_j(x) = \sigma \left( \frac{x - \mu_j}{s} \right) \quad \sigma_a = \frac{1}{1 + \exp(-a)}$$

# Least-squares

- Loss function

$$L(t, y(x)) = \{y(x) - t\}^2 \text{ (squared residuals)}$$

- Empirical risk

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2$$

# Least-squares

- Solution

$$\mathbf{w}^{\star} = (\mathbf{\Phi}^T \mathbf{\Phi})^{-1} \mathbf{\Phi}^T \mathbf{t}$$

$$\mathbf{\Phi} = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

# Least-squares and Maximum Likelihood

- Review for maximum-likelihood estimation (MLE)

$$q(\mathbf{x}; \boldsymbol{\theta}) \longrightarrow p(\mathbf{x})$$

$$\mathcal{D} = \{\mathbf{x}\}_{i=1}^n$$

$$\hat{\boldsymbol{\theta}}_{\text{ML}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(\mathcal{D} \mid \boldsymbol{\theta}) = \prod_{i=1}^n q(\mathbf{x}; \boldsymbol{\theta})$$

$$\hat{\boldsymbol{\theta}}_{\text{ML}} = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmax}} \log L(\boldsymbol{\theta}) = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmax}} \left[ \sum_{i=1}^n \log q(\mathbf{x}_i; \boldsymbol{\theta}) \right]$$

# Least-squares and Maximum Likelihood

- Relation between LS and MLE

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon \quad \epsilon \sim \mathcal{N}(0, \beta^{-1})$$

$$p(t \mid \mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t \mid y(\mathbf{x}, \mathbf{w}), \beta^{-1})$$

$$\mathbb{E}[t \mid \mathbf{x}] = \int t p(t \mid \mathbf{x}) dt = y(\mathbf{x}, \mathbf{w})$$

$$p(\mathbf{t} \mid \mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n \mid \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1})$$

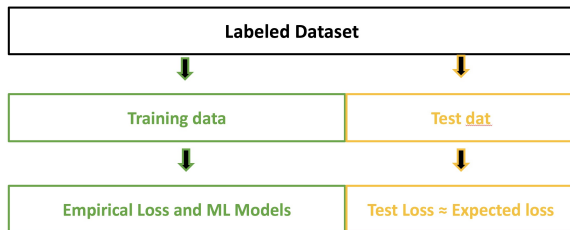
$$\begin{aligned} \ln p(\mathbf{t} \mid \mathbf{w}, \beta) &= \sum_{n=1}^N \ln \mathcal{N}(t_n \mid \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}) \\ &= \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(\mathbf{w}) \end{aligned}$$

$$\mathbf{w}_{MLE}^* = \arg \max_{\mathbf{w}} \ln p(\mathbf{t} \mid \mathbf{w}, \beta)$$

$$\mathbf{w}_{MLE}^* = \mathbf{w}_{LS}^*$$



# Model Evaluation



- When learning a model, you should pretend that you don't have the test data yet. If the test-set labels influence the learned model in any way, accuracy estimates will be biased.
- Your test set should be large enough to detect meaningful changes in the accuracy of your algorithm, but not necessarily much larger.
- When randomly selecting training or validation sets, we may want to ensure that class proportions are maintained in each selected set.

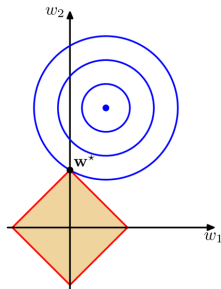
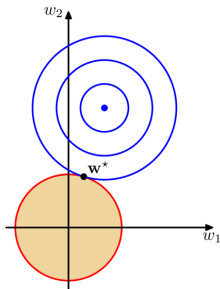
## Ridge regression

$$\begin{aligned} & E_D(\mathbf{w}) + \lambda E_W(\mathbf{w}) \\ &= \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \\ \mathbf{w}_{\text{ridge}}^* &= (\lambda \mathbf{I} + \boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^T \mathbf{t} \end{aligned}$$

# LASSO

$$E_D(\mathbf{w}) + \lambda E_W(\mathbf{w})$$

$$= \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \Phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|^q$$



## RR and Maximum a Posteriori Estimation (MAP)

- Review for maximum a posteriori estimation (MAP)

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

Prior  $p(\boldsymbol{\theta})$

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

$$\hat{\boldsymbol{\theta}}_{\text{MAP}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(\boldsymbol{\theta}|\mathcal{D})$$

$$\hat{\boldsymbol{\theta}}_{\text{MAP}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \left( \sum_{i=1}^n \log q(\mathbf{x}_i|\boldsymbol{\theta}) + \log p(\boldsymbol{\theta}) \right)$$

## RR and Maximum a Posteriori Estimation (MAP)

- Relation between RR and MAP

$$p(\mathbf{t} \mid \mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n \mid \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1})$$

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} \mid \mathbf{m}_0, \mathbf{S}_0)$$

$$p(\mathbf{w} \mid \mathbf{t}) = \mathcal{N}(\mathbf{w} \mid \mathbf{m}_N, \mathbf{S}_N)$$

$$\mathbf{m}_N = \mathbf{S}_N (\mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \boldsymbol{\Phi}^T \mathbf{t}), \quad \mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}$$

# RR and Maximum a Posteriori Estimation (MAP)

- Relation between RR and MAP

$$p(\mathbf{t} \mid \mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n \mid \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1})$$

$$p(\mathbf{w} \mid \alpha) = \mathcal{N}(\mathbf{w} \mid \mathbf{0}, \alpha^{-1} \mathbf{I})$$

$$p(\mathbf{w} \mid \mathbf{t}) = \mathcal{N}(\mathbf{w} \mid \mathbf{m}_N, \mathbf{S}_N)$$

$$\mathbf{m}_N = \beta \mathbf{S}_N \boldsymbol{\Phi}^T \mathbf{t}, \quad \mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}$$

$$\mathbf{w}_{MAP}^* = \arg \max_{\mathbf{w}} \ln p(\mathbf{w} \mid \mathbf{t})$$

$$= \arg \max_{\mathbf{w}} \left( -\frac{\beta}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n)\}^2 - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + \text{const} \right)$$

$$\mathbf{w}_{MAP}^* = \mathbf{w}_{ridge}^*$$

# Bayesian Models

- Bayesian Estimation

$$\int \boldsymbol{\theta} p(\boldsymbol{\theta}|\mathcal{D}) d\boldsymbol{\theta} \quad (\text{posterior expectation})$$

- Bayesian Estimation

$$\begin{aligned}\hat{p}_{\text{Bayes}}(\mathbf{x}) &= \int q(\mathbf{x}|\boldsymbol{\theta}) p(\boldsymbol{\theta}|\mathcal{D}) d\boldsymbol{\theta} \\ &= \int q(\mathbf{x}|\boldsymbol{\theta}) \frac{p(\mathcal{D}|\boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\mathcal{D})} d\boldsymbol{\theta} = \int q(\mathbf{x}|\boldsymbol{\theta}) \frac{\prod_{i=1}^n q(\mathbf{x}_i|\boldsymbol{\theta}) p(\boldsymbol{\theta})}{\int \prod_{i=1}^n q(\mathbf{x}_i|\boldsymbol{\theta}') p(\boldsymbol{\theta}') d\boldsymbol{\theta}'} d\boldsymbol{\theta}\end{aligned}$$

# Bayesian Linear Regression

$$p(t \mid \mathbf{t}, \alpha, \beta) = \int p(t \mid \mathbf{w}, \beta) p(\mathbf{w} \mid \mathbf{t}, \alpha, \beta) d\mathbf{w} \quad (\text{Predictive distribution})$$

$$p(t \mid \mathbf{x}, \mathbf{t}, \alpha, \beta) = \mathcal{N}(t \mid \mathbf{m}_N^T \phi(\mathbf{x}), \sigma_N^2(\mathbf{x}))$$

$$\sigma_N^2(\mathbf{x}) = \frac{1}{\beta} + \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x})$$

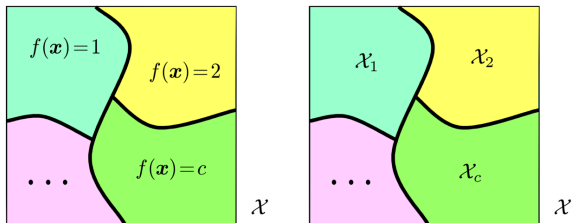


# Outline

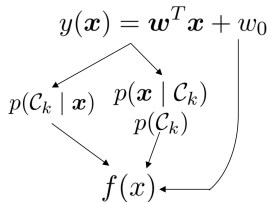
Linear Regression

**Linear Classification**

# What is Linear classification



- Probabilistic Discriminative Models
- Probabilistic Generative Models
- Discriminant Functions



## Least squares for classification?

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}$$

$$\mathbf{y}(\mathbf{x}) = \tilde{\mathbf{W}}^T \tilde{\mathbf{x}}, \quad \tilde{\mathbf{w}}_k = (w_{k0}, \mathbf{w}_k^T)^T, \quad \tilde{\mathbf{x}} = (1, \mathbf{x}^T)^T$$

$$\{\mathbf{x}_n, \mathbf{t}_n\}, n = 1, \dots, N$$

$$\tilde{\mathbf{X}} - n^{\text{th}} \text{ row} - \tilde{\mathbf{x}}_n^T$$

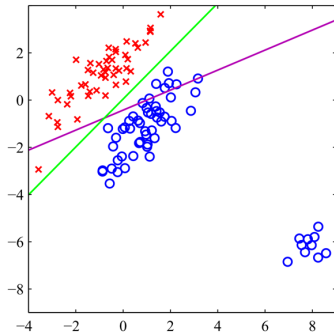
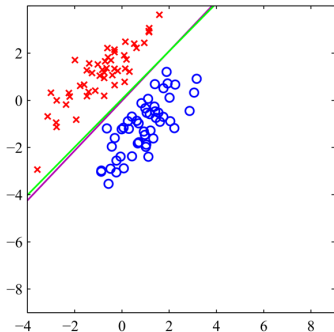
$$\mathbf{T} - n^{\text{th}} \text{ row} - \mathbf{t}_n^T$$

$$E_D(\tilde{\mathbf{W}}) = \frac{1}{2} \text{Tr} \left\{ (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T})^T (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T}) \right\}$$

$$\tilde{\mathbf{W}} = \left( \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}^T \mathbf{T} = \tilde{\mathbf{X}}^\dagger \mathbf{T}$$

$$y(\mathbf{x}) = \tilde{\mathbf{W}}^T \tilde{\mathbf{x}} = \mathbf{T}^T \left( \tilde{\mathbf{X}}^\dagger \right)^T \tilde{\mathbf{x}}$$

## Least squares for classification?



least squares is highly sensitive to outliers

# Probabilistic Discriminative Models

- Logistic regression

$$p(\mathcal{C}_1 | \phi) = \sigma(\mathbf{w}^T \phi)$$

$$\sigma(a) = \frac{1}{1 + \exp(-a)}$$

$$p(\mathcal{C}_2 | \phi) = 1 - p(\mathcal{C}_1 | \phi)$$

Why sigmoid function?

$$p(\mathcal{C}_1 | \mathbf{x}) = \frac{p(\mathbf{x} | \mathcal{C}_1) p(\mathcal{C}_1)}{p(\mathbf{x} | \mathcal{C}_1) p(\mathcal{C}_1) + p(\mathbf{x} | \mathcal{C}_2) p(\mathcal{C}_2)}$$

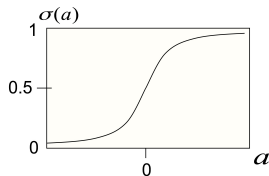
$$p(\mathcal{C}_1 | \mathbf{x}) = \frac{1}{1 + \exp(-a)} = \sigma(a)$$

$$a = \ln \frac{p(\mathbf{x} | \mathcal{C}_1) p(\mathcal{C}_1)}{p(\mathbf{x} | \mathcal{C}_2) p(\mathcal{C}_2)} = \ln \frac{p(\mathcal{C}_1 | \mathbf{x})}{p(\mathcal{C}_2 | \mathbf{x})}$$

$$p(\mathcal{C}_k | \mathbf{x}) \sim \text{Ber}(\sigma(\mathbf{w}^T \phi))$$

$\downarrow$

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



# Probabilistic Discriminative Models

- Logistic regression

$$p(\mathbf{t} \mid \mathbf{w}) = \prod_{n=1}^N \{p(\mathcal{C}_1 \mid \phi_n)\}^{t_n} \{1 - p(\mathcal{C}_1 \mid \phi_n)\}^{1-t_n}$$

$$y_n = p(\mathcal{C}_1 \mid \phi_n)$$

$$E(\mathbf{w}) = -\ln p(\mathbf{t} \mid \mathbf{w}) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln (1 - y_n)\}$$

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n$$

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n$$

# Probabilistic Generative Models

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

↓

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

$$p(\mathcal{C}_k)$$

↓

$$f(x)$$

# Probabilistic Generative Models

- Linear discriminant

$$p(\mathbf{x} \mid \mathcal{C}_k) = \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\}$$

(assuming that features are continuous and all classes share the same covariance matrix)



# Probabilistic Generative Models

- Linear discriminant

$$p(\mathbf{x} | \mathcal{C}_k) = \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\}$$

(assuming that features are continuous and all classes share the same covariance matrix)

Linear?

$$p(\mathcal{C}_1 | \mathbf{x}) = \frac{p(\mathbf{x} | \mathcal{C}_1) p(\mathcal{C}_1)}{p(\mathbf{x} | \mathcal{C}_1) p(\mathcal{C}_1) + p(\mathbf{x} | \mathcal{C}_2) p(\mathcal{C}_2)} = \frac{1}{1 + \exp(-a)} = \sigma(a)$$

$$a = \ln \frac{p(\mathbf{x} | \mathcal{C}_1) p(\mathcal{C}_1)}{p(\mathbf{x} | \mathcal{C}_2) p(\mathcal{C}_2)}$$

$$p(\mathcal{C}_1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

$$\mathbf{w} = \Sigma^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \quad w_0 = -\frac{1}{2} \boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_2^T \Sigma^{-1} \boldsymbol{\mu}_2 + \ln \frac{p(\mathcal{C}_1)}{p(\mathcal{C}_2)}$$

# Probabilistic Generative Models

- Linear discriminant

$$p(\mathbf{x} | \mathcal{C}_k) = \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\}$$

(assuming that features are continuous and all classes share the same covariance matrix)

Linear?

$$p(\mathcal{C}_k | \mathbf{x}) = \frac{p(\mathbf{x} | \mathcal{C}_k) p(\mathcal{C}_k)}{\sum_j p(\mathbf{x} | \mathcal{C}_j) p(\mathcal{C}_j)} = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

$$a_k = \ln p(\mathbf{x} | \mathcal{C}_k) p(\mathcal{C}_k)$$

$$a_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}$$

$$\mathbf{w}_k = \Sigma^{-1} \boldsymbol{\mu}_k \quad w_{k0} = -\frac{1}{2} \boldsymbol{\mu}_k^T \Sigma^{-1} \boldsymbol{\mu}_k + \ln p(\mathcal{C}_k)$$

# Probabilistic Generative Models

- Maximum likelihood solution for Linear discriminant

$$\{\mathbf{x}_n, t_n\}_{n=1}^N, t_n = 1 \longleftrightarrow \mathcal{C}_1, t_n = 0 \longleftrightarrow \mathcal{C}_2$$

$$p(\mathcal{C}_1) = \pi, p(\mathcal{C}_2) = 1 - \pi$$

$$p(\mathbf{x}_n, \mathcal{C}_1) = p(\mathcal{C}_1) p(\mathbf{x}_n | \mathcal{C}_1) = \pi \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma})$$

$$p(\mathbf{x}_n, \mathcal{C}_2) = p(\mathcal{C}_2) p(\mathbf{x}_n | \mathcal{C}_2) = (1 - \pi) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma})$$

$$p(\mathbf{t}, \mathbf{X} | \pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) = \prod_{n=1}^N [\pi \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma})]^{t_n} [(1 - \pi) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma})]^{1-t_n}$$

$$\pi = \frac{1}{N} \sum_{n=1}^N t_n = \frac{N_1}{N} = \frac{N_1}{N_1 + N_2}, \boldsymbol{\mu}_1 = \frac{1}{N_1} \sum_{n=1}^N t_n \mathbf{x}_n, \boldsymbol{\mu}_2 = \frac{1}{N_2} \sum_{n=1}^N (1 - t_n) \mathbf{x}_n$$

$$\mathbf{S} = \frac{N_1}{N} \mathbf{S}_1 + \frac{N_2}{N} \mathbf{S}_2$$

$$\mathbf{S}_1 = \frac{1}{N_1} \sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \boldsymbol{\mu}_1) (\mathbf{x}_n - \boldsymbol{\mu}_1)^T, \mathbf{S}_2 = \frac{1}{N_2} \sum_{n \in \mathcal{C}_2} (\mathbf{x}_n - \boldsymbol{\mu}_2) (\mathbf{x}_n - \boldsymbol{\mu}_2)^T$$

# Probabilistic Generative Models

- Naïve-Bayes (NB) classifier
- Naïve-Bayes (NB) assumption

$$x_i \perp x_{\{j \neq i\}} \mid t$$

- Bernoulli NB classifier

$$x_i \in \{0, 1\} \ \& \ p(x_i \mid \mathcal{C}_k) \sim \text{Ber}(\mu_{ki})$$

$$p(\mathbf{x} \mid \mathcal{C}_k) = \prod_{i=1}^D \mu_{ki}^{x_i} (1 - \mu_{ki})^{1-x_i}$$

$$a_k = \ln p((\mathbf{x} \mid \mathcal{C}_k) p(\mathcal{C}_k))$$

$$a_k(\mathbf{x}) = \sum_{i=1}^D \{x_i \ln \mu_{ki} + (1 - x_i) \ln (1 - \mu_{ki})\} + \ln p(\mathcal{C}_k)$$

# Hinge Loss and Support Vector Machines

- Loss Functions for Classification

$$t_n \in \{-1, 1\}$$

$$y_n > 0 \quad \hat{t}_n = 1, y_n < 0 \quad \hat{t}_n = -1$$

- 0-1 loss

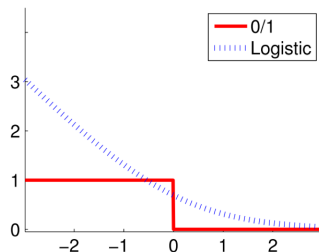
$$E_{0/1}(t_n, y_n) = 1 - \text{sign}\{t_n y(\mathbf{x}_n)\}$$

- Log loss

$$E_{\log}(t_n, y_n) = \ln\{1 + \exp(-y_n t_n)\}$$

equals to

$$E_{\text{cross-ent}}(t_n, y_n) = \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\} \quad (t_n \in \{0, 1\})$$



# Hinge Loss and Support Vector Machines

- Loss Functions for Classification

- Hinge Loss

$$t_n \in \{-1, 1\}$$

$$y_n > 0 \rightarrow \hat{t}_n = 1, y_n < 0 \rightarrow \hat{t}_n = -1$$

$$E_{\text{Hinge}}(t_n, y_n) = [1 - y_n t_n]_+$$

$[\cdot]_+$  denotes the positive part

- Support Vector Classifier

$$L_{\text{SVC}} = \sum_{n=1}^N E_{\text{Hinge}}(t_n, y_n) + \lambda \|\mathbf{w}\|^2$$

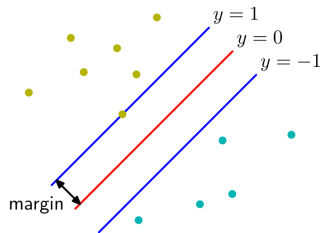
# Hinge Loss and Support Vector Machines

- Maximum-Margin View for SVC

$$\arg \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n [t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b)] \right\}$$

s.t.  $t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 0, n = 1, \dots, N$

$$\text{s.t. } \min_n [t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b)] = 1$$

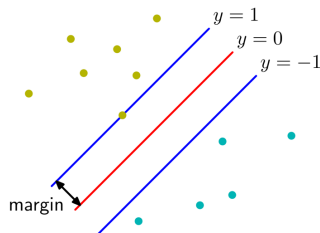


# Hinge Loss and Support Vector Machines

- Maximum-Margin View for SVC

$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

$$\text{s.t. } t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1, n = 1, \dots, N$$





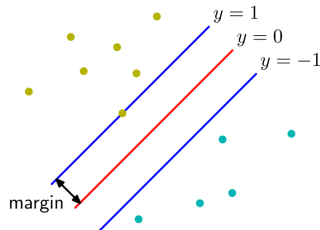
# Hinge Loss and Support Vector Machines

- Maximum-Margin View for SVC

$$\arg \min_{\mathbf{w}, b, \xi} C \sum_{n=1}^N \xi_n + \frac{1}{2} \|\mathbf{w}\|^2$$

$$\text{s.t. } t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1 - \xi_n, n = 1, \dots, N$$

$$\xi_n \geq 0$$



# Model Evaluation for Classification

- Performance Matrices

- Confusion matrix

		Actual	
		Class +	Class -
Predicted	Class +	TP	FP
	Class -	FN	TN

- Accuracy

$$\frac{TP + TN}{TP + FP + FN + TN}$$

- Error rate

$$\frac{FP + FN}{TP + FP + FN + TN}$$

# Model Evaluation for Classification

- Performance Matrices

- Confusion matrix

		Actual	
		Class +	Class -
Predicted	Class +	TP	FP
	Class -	FN	TN

- Precision

$$\text{TP} / (\text{TP} + \text{FP})$$

- Recall

$$\text{TP} / (\text{TP} + \text{FN})$$

- F-measure

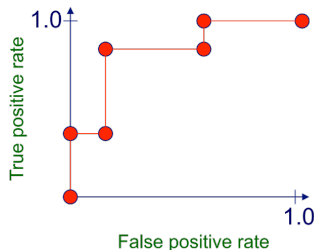
$$F_{\beta} = (1 + \beta^2) \cdot \frac{\text{precision} \cdot \text{recall}}{(\beta^2 \cdot \text{precision}) + \text{recall}}$$

# Model Evaluation for Classification

- Performance Matrices

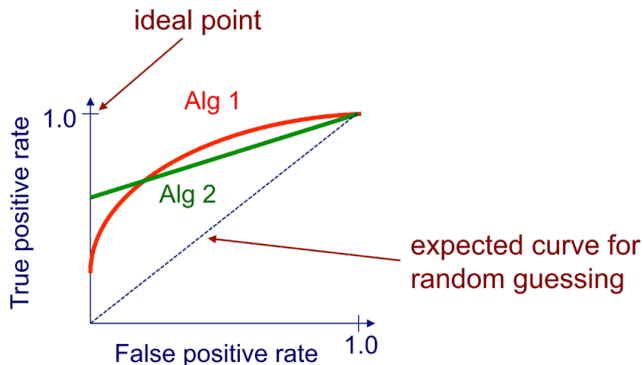
- ROC

instance	confidence positive	correct class
Ex 9	.99	+
Ex 7	.98	+
Ex 1	.72	-
Ex 2	.70	+
Ex 6	.65	+
Ex 10	.51	-
Ex 3	.39	-
Ex 5	.24	+
Ex 4	.11	-
Ex 8	.01	-



# Model Evaluation for Classification

- Performance Matrices
  - ROC



# Model Evaluation for Classification

- Performance Matrices

- AUC

