

教材: 筒老师手稿 + ppt. (无教材, 参考教材 PRML)
(Bishop' 2006)

Day 1:

1. Introduction
2. Probability Distribution (data & compute)
3. Linear Regression (输出连续)
4. Linear Classification (输出离散)
5. Kernel Method. (another trick)
6. Sparse Kernel Method \equiv SVM \rightarrow support vector
7. Mixture Model & EM.
8. Approximate Inference.

1. Introduction



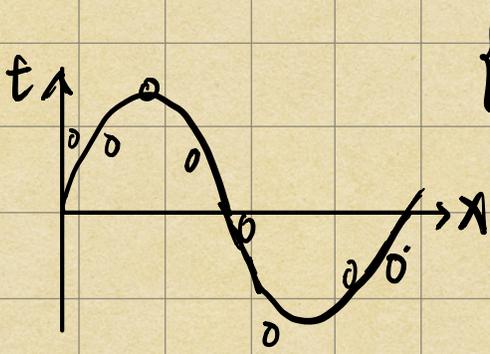
Unsupervised \rightarrow cluster
without teacher

$$\{(x_1, t_1), (x_2, t_2), \dots, (x_n, t_n)\} \Rightarrow X \xrightarrow{\text{mapping } f} t$$

1' Random variable

2' train / testing \rightarrow {prediction / generalization (泛化)}

2. Generalization: fitting $\sin(x)$



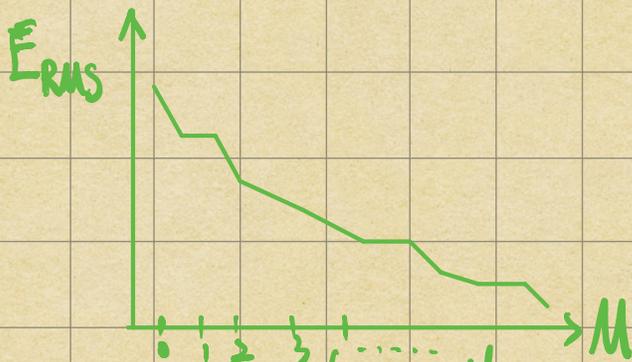
$$\{X_n, t_n\}^N \quad y(x, w) \rightarrow t$$

1. Assume $y(x, w) = w_0 + w_1 x + \dots + w_m x^m$
2. Learning Criterion / Objective: Error Function & Divergence.

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

root-mean-squares:

$$\sqrt{2 \cdot E(w^*) / N} = E_{RMS}$$



\Rightarrow larger M , more complex the Model, less Error (Overfitting)

how to trade-off the overfitting and underestimating?

$$\tilde{E} = \frac{1}{2} \sum_{n=1}^N \{y(x_n, w) - t_n\}^2 + \frac{\lambda}{2} \|w\|^2 \rightarrow \text{adjust error and } w.$$

hyper parameter. (regularization term) = λ, M

$$\|w\| = W^T \cdot W$$

regularization

reduce the value of w when M increase

Add probability into it:

Prior / likelihood / Posterior

$$P(w|D) = \frac{P(D|w) \cdot P(w)}{P(D)} \rightarrow \text{in disperse condition: } P(D) = \int P(D|w) \cdot P(w) \cdot P(D|w) dw$$

(evidence)

$$P(w|D) \propto P(D|w) \cdot P(w) \quad \text{For model selection}$$

Supplement: $N(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \cdot \exp(-\frac{1}{2\sigma^2}(x-\mu)^2) > 0$

$$y(x, w) \quad (-R^p \quad -R^{M+1})$$

3. with D-dimension Gaussian:

$$N(x|\mu, \Sigma) = \frac{1}{(2\pi)^{\frac{D}{2}}} \cdot \frac{1}{|\Sigma|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right\}$$

with 1-dimension, we obtain probability of dataset:

$$p(x|\mu, \sigma^2) = \prod_{n=1}^N N(x_n|\mu, \sigma^2)$$

then likelihood:

$$\begin{aligned} \ln P(x|\mu, \sigma^2) &= \ln(N(x_1|\mu, \sigma^2) \cdot N(x_2|\mu, \sigma^2) \cdots N(x_n|\mu, \sigma^2)) \\ &= -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi) \end{aligned}$$

Maximum Likelihood:

$$\frac{d}{d\mu} (\ln P(x|\mu, \sigma^2)) = \frac{1}{\sigma^2} \sum_{n=1}^N (x_n - \mu) \xrightarrow{\text{let } = 0} \sum_{n=1}^N x_n = \sum_{n=1}^N \mu$$

$$\therefore \mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n, \text{ likewise } \sigma_{ML}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})^2$$

★ **Maximum Likelihood** (especially W)

$$\mu^* = \frac{1}{N} \cdot \sum_{n=1}^N x_n = \mu$$

$$\text{Define: } \Sigma^* = \frac{1}{N-1} \sum_{n=1}^N (x_n - \mu^*)^2$$

$$P(t|x, W, \beta) = N(t_n | y(x_n, W), \beta^{-1}) \quad \beta^{-1} = \frac{1}{\sigma^2} = \frac{1}{2} \text{ as precision}$$

$$\therefore P(D|W) = P(t|x, W, \beta) = \prod_{n=1}^N N(t_n | x_n, W, \beta)$$

$$\therefore \ln P(t|x, W, \beta) = -\frac{\beta}{2} \sum_{n=1}^N (y(x_n, W) - t_n)^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$$

$$\text{Let } \beta^{-1} = \frac{1}{N} \sum_{n=1}^N \{y(x_n, W_{ML}) - t_n\}^2$$

$$\theta E(W) \leftrightarrow \ln P(D|W)$$

$$W_{MAP} = \underset{W}{\text{argmax}} P(W|x, t, \alpha, \beta) = \underset{W}{\text{argmin}} \left\{ \frac{\beta}{2} \sum_{n=1}^N (y(x_n, W) - t_n)^2 + \frac{\alpha}{2} \cdot W^T \cdot W \right\}$$

Maximum a Posterior

4. Proof above: Given $\ln p(t|x, \omega, \beta)$
 $= -\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \omega) - t_n\}^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$
 β is irrelevant to ω_{ML} , after we calculate ω_{ML}
then $\beta_{ML} = \frac{1}{N} \sum_{n=1}^N \{y(x_n, \omega_{ML}) - t_n\}^2$
finally: we obtain "predictive distribution"
 $p(t|x, \omega_{ML}, \beta_{ML}) = N(t|y(x, \omega_{ML}), \beta_{ML}^{-1})$

Section MAP:

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

$\ln(p(w|\alpha)) = \frac{M+1}{2} \ln \frac{\alpha}{2\pi} - \frac{\alpha}{2} w^T \cdot w$

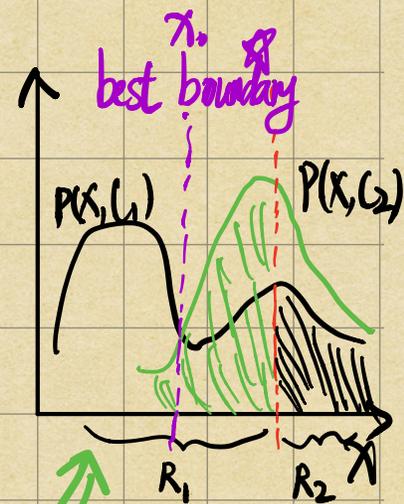
$\therefore P(w|x, t, \alpha, \beta) \propto P(t|x, w, \beta) P(w|\alpha)$
↓ posterior distrib ↓ prior distrib ↓ likelihood

Data flow: $w_{ML} \xrightarrow{P(t|x)} \beta(\text{precision}) \xrightarrow{P(w|\alpha)} w_{MAP}$

Decision Theory:

C (presence of cancer) $\begin{cases} C_1: \text{with cancer} \\ C_2: \text{without cancer} \end{cases}$

$\hat{C}_{MAP} = \arg \max_k P(C_k|x) = \arg \max_k \frac{P(C_k) \cdot P(x|C_k)}{P(x)}$



1. Minimizing the misclassification Rate

$P(\text{mistake}) = P(x \in R_1, C_2) + P(x \in R_2, C_1)$
 $= \int_{R_1} P(x, C_2) dx + \int_{R_2} P(x, C_1) dx$

★ Minimizing the expected loss

good. Loss: Minimum Bayes Risk \rightarrow $\langle MBR \rangle$
Expected Loss
cancer normal.

5.

Loss matrix cancer normal $\begin{pmatrix} 0 & 1000 \\ 1 & 0 \end{pmatrix} = L = [L_{kj}]$

$$E[L] = \sum_k \sum_j \int_{R_j} L_{kj} \cdot P(x, C_k) \cdot dx \Rightarrow \hat{C}_j = \underset{P(x) \cdot P(C_k|x)}{\text{argmin}} \sum_k L_{kj} P(C_k|x)$$

$P(C_k|x) = \frac{P(x|C_k) \cdot P(C_k)}{P(x)}$ → generative model (indirect model)

↓
discriminative model (direct model)

* loss function for regression:

$$E[L] = \iint L(t, y(x)) \cdot P(x, t) dx dt$$

$$= \iint \{y(x) - t\}^2 \cdot P(x, t) dx dt \quad \int y(x) P(x, t) dt = \int t \cdot P(x, t) dt$$

$$\frac{\partial E[L]}{\partial y(x)} = 2 \int \{y(x) - t\} P(x, t) dt = 0$$

$$y^*(x) = \frac{\int t P(x, t) dt}{P(x)} = E_t[t|x]$$



non-prob
LS (least square)
↓ upgrade
RLS (regularization
least square)

probability
ML (Maximum likelihood)
↓
MAP (Maximum a posterior)

Supplement:

$$P(t|x, X, t) = \int P(t|x, w) \cdot P(w|x, t) dw$$

test data training data = $N(t | m(x), s^2(x))$

$m(x) =$
 $s^2(x) =$ ★

Information Theory: "Entropy" \leftrightarrow "uncertainty"

KL "divergence":
Kullback-Leibler

Definition: Two unrelated events, $h(x,y) = h(x) + h(y)$ & $h(x) = -\log_2 P(x)$ ↓ bits

$$H[X] = -\sum_x P(x) \cdot \log_2 P(x) \quad [H \text{ is a expectation of information}]$$

$$H[P] = -\sum_i P(x_i) \cdot \ln P(x_i) = -\int p(x) \cdot \ln p(x) \cdot dx$$

$$= E[-\ln P(x)]$$

for P is \ln
for x_i is \log_2

Maximum Entropy (ME)

$$\tilde{H} = -\sum P(x_i) \ln P(x_i) + \lambda (\sum P(x_i) - 1)$$

Given: X is continuous random variables:

$$\Rightarrow H[X] = -\int p(x) \cdot \ln p(x) \cdot dx \quad (\text{differential entropy})$$

Given three constraints we have:

$$-\int p(x) \ln p(x) dx + \lambda_1 (\int_{-\infty}^{+\infty} p(x) dx - 1) + \lambda_2 (\int_{-\infty}^{+\infty} x p(x) dx - \mu) + \lambda_3 (\int_{-\infty}^{+\infty} (x-\mu)^2 p(x) dx - \sigma^2)$$

$$\Rightarrow P(x) = \exp \{ -1 + \lambda_1 + \lambda_2 x + \lambda_3 (x-\mu)^2 \}$$

$$\text{finally we obtain } \Rightarrow P(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left(-\frac{(x-\mu)^2}{2\sigma^2} \right) \quad \times$$

$$\text{then } H[X] = \frac{1}{2} \{ 1 + \ln(2\pi\sigma^2) \}$$

Entropy versus Pattern Recognition:

$$KL(p \parallel q) = -\int p(x) \ln q(x) dx - (-\int p(x) \ln p(x) dx)$$

$$= -\int p(x) \ln \frac{q(x)}{p(x)} dx \geq 0$$

unknown approx true distrib approx distrib $\int p(x) \ln p(x) dx = 0$

$KL(p||q) \neq KL(q||p)$

proof: with Jensen's Inequality we obtain:

$f(\lambda a + (1-\lambda)b) \leq \lambda f(a) + (1-\lambda)f(b)$

$\Rightarrow f(E[X]) \leq E[f(X)]$

$KL(p||q) = - \int p(x) \cdot \ln \left\{ \frac{q(x)}{p(x)} \right\} dx \geq - \ln \int q(x) \cdot dx = 0$

Parametric Model: estimate θ

$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} KL(p(x)|q(x|\theta)) \quad \& \quad KL(p||q) \approx \sum_{i=1}^N \{-\ln q(x_i|\theta) + \ln p(x_i)\}$

Mutual Information: according to independence

$I(x,y) = KL(p(x,y) || p(x) \cdot p(y))$
 $= - \iint p(x,y) \ln \left(\frac{p(x) \cdot p(y)}{p(x,y)} \right) \geq 0$

Day 2: Probability Distribution:

① Binomial - beta

Intro: Discrete: Natural Language

② Multinomial - Dirichlet

likelihood-prior conjugate-prior

Continuous:

③ Gaussian - Gaussian - Gamma

sequential learning

④ Student's t



Exponential Family

Background:

$x = [x_a^T \ x_b^T]^T \left(\begin{bmatrix} x_a \\ x_b \end{bmatrix} \right) \propto \frac{P(x_a|x_a)}{P(x_a|x_b)}$

1. Bernoulli Distribution $X \in \{0, 1\}$

Method 1: Maximum Likelihood

In Bernoulli, Given Dataset $D = \{X_1, \dots, X_N\}$, we obtain the likelihood function $P(D|\mu) = \prod_{n=1}^N P(X_n|\mu) = \prod_{n=1}^N \mu^{X_n} (1-\mu)^{1-X_n}$

$$\Rightarrow \ln P(D|\mu) = \sum_{n=1}^N \{X_n \ln \mu + (1-X_n) \ln(1-\mu)\} \stackrel{\text{let}}{=} 0$$

$$\therefore \mu_{ML} = \frac{1}{N} \sum_{n=1}^N X_n \quad \text{defection: overfitting while given a small dataset.}$$

Method 2: Posterior Bayesian.

$$\text{Bin}(m|N, \mu) = \binom{N}{m} \mu^m (1-\mu)^{N-m} \quad \text{where } \binom{N}{m} = \frac{N!}{(N-m)!m!}$$

let introduce prior probability distrib $P(\mu)$:

$P(\mu)$ must have conjugacy for the same functional formation!

Definition:

$$\text{Beta}(\mu|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} \cdot (1-\mu)^{b-1}$$

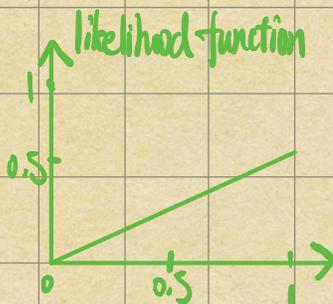
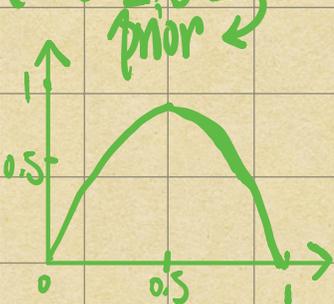
$$E[\mu] = \frac{a}{a+b} \quad \text{hyperparam} \quad \text{var}[\mu] = \frac{ab}{(a+b)^2(a+b+1)}$$

$$P(\mu|m, l, a, b) \propto \text{Bin}(m|N, \mu) \cdot \text{Beta}(\mu|a, b)$$

$$P(\mu|m, l, a, b) \propto \mu^{m+a-1} \cdot (1-\mu)^{l+b-1} \quad (\text{without regularization})$$

$$\therefore P(\mu|m, l, a, b) \propto \frac{\Gamma(m+a)\Gamma(l+b)}{\Gamma(m+a+l+b)} \mu^{m+a-1} \cdot (1-\mu)^{l+b-1}$$

with: $a=2, b=2$
prior



Conjugate Prior:

1. Sequential Learning: Update.

2. Prediction: $\int_0^1 \dots \int_0^1 dw$ integrand.

Example:

$$\lim_{m, l \rightarrow \infty} P(X=1 | D) = \int_0^1 P(X=1 | \mu) \cdot P(\mu | D) \cdot d\mu = \int_0^1 \mu \cdot P(\mu | D) d\mu$$
$$= E[\mu | D] = \frac{m+a}{m+l+1} \rightarrow \frac{m}{N}$$

↑
predictive distrib

2. Multinomial Variables:

$$\text{Mult}(m_1, m_2, \dots, m_k | \mu, N) = \binom{N}{m_1, m_2, \dots, m_k} \prod_{k=1}^k \mu_k^{m_k}, \text{ where}$$

$$\binom{N}{m_1, m_2, \dots, m_k} = \frac{N!}{m_1! m_2! \dots m_k!} \text{ and } \sum_{k=1}^k m_k = N$$

Dirichlet distribution: (prior for mult variables)

$$P(\mu | \alpha) \propto \prod_{k=1}^k \mu_k^{\alpha_k - 1} \quad 0 \leq \mu_k \leq 1 \quad \sum_k \mu_k = 1$$

to simplex: \rightarrow (a bounded linear manifold)

$$\text{Dir}(\mu | \alpha) = \frac{P(\alpha_1 + \dots + \alpha_k)}{P(\alpha_1) \dots P(\alpha_k)} \prod_{k=1}^k \mu_k^{\alpha_k - 1}$$

Normalization Term.

$$P(\mu | D, \alpha) \propto P(D | \mu) \cdot P(\mu | \alpha) \propto \prod_{k=1}^k \mu_k^{\alpha_k + m_k - 1}$$

$$\Rightarrow P(\mu | D, \alpha) = \text{Dir}(\mu | \alpha + m) = \frac{P(\alpha_0 + N)}{P(\alpha_1 + m) \dots P(\alpha_k + m_k)} \prod_{k=1}^k \mu_k^{\alpha_k + m_k - 1}$$

3. Gaussian Distribution:

$$N(x|\mu, \Sigma) = \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x-\mu)^T \cdot \Sigma^{-1} (x-\mu) \right\}$$

where Σ is a $D \times D$ variance matrix $|\Sigma|$ is determinant

Δ : Mahalanobis distance
(二次型形式)

For Eigenvector Equation:

$$\sum U_i = \lambda_i U_i$$

symmetric eigenvector eigenvalue

from general to isotropic



Conditional Gaussian Distribution: $P(x_a|x_b)$

Given: $N(x|\mu, \Sigma)$, $x = \begin{pmatrix} x_a \\ x_b \end{pmatrix}$, $\mu = \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix}$, $\Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}$

Let $\Lambda = \Sigma^{-1}$ (precision matrix)

$$\begin{aligned} -\frac{1}{2} (x-\mu)^T \cdot \Sigma^{-1} \cdot (x-\mu) &= -\frac{1}{2} (x_a - \mu_a)^T \Lambda_{aa} (x_a - \mu_a) \\ &\quad -\frac{1}{2} (x_a - \mu_a)^T \Lambda_{ab} (x_b - \mu_b) \\ &\quad -\frac{1}{2} (x_b - \mu_b)^T \Lambda_{ba} (x_a - \mu_a) \\ &\quad -\frac{1}{2} (x_b - \mu_b)^T \Lambda_{bb} (x_b - \mu_b) \end{aligned}$$

After complicative equation, we obtain:

$$\mu_{ab} = \mu_a + \Sigma_{aa} \Sigma_{bb}^{-1} (x_b - \mu_b)$$

$$\Sigma_{ab} = \Sigma_{aa} - \Sigma_{aa} \Sigma_{bb}^{-1} \Sigma_{ba}$$

Every quadratic equation could be expressed by a form of Gaussian Distribution via x^2 for Σ , x' for mean.

Marginal Gaussian Distribution:

$$P(x_a) = \int P(x_a, x_b) dx_b$$

We are integrating out x_b , we find terms involve x_b :

$$-\frac{1}{2} x_b^T \Lambda_{bb} x_b + x_b^T \underbrace{\{ \Lambda_{bb} \mu_b - \Lambda_{ba} (x_a - \mu_a) \}}_m$$

$$= -\frac{1}{2} (x_b - \Lambda_{bb}^{-1} m)^T \Lambda_{bb} (x_b - \Lambda_{bb}^{-1} m) + \frac{1}{2} m^T \Lambda_{bb}^{-1} m$$

$$\Rightarrow \int \exp \left\{ -\frac{1}{2} (x_b - \Lambda_{bb}^{-1} m)^T \Lambda_{bb} (x_b - \Lambda_{bb}^{-1} m) \right\} dx_b$$

then picking up terms depending on x_a :

$$\Rightarrow P(x_a) = N(x_a | \mu_a, \Sigma_a)$$

★★ Sequential estimation:

Sequential methods allow data points to be processed one at a time and then discarded \Rightarrow on-line application.

Given Conjugate Prior: $P(\mu) = N(\mu | \mu_0, \sigma_0^2)$

Bayesian inference for Gaussian:

Case 1: unknown mean, known variance (σ^2)

$$\Rightarrow P(X | \mu) = \prod_{n=1}^N P(x_n | \mu) = \frac{1}{\sqrt{2\pi} \sigma} \exp \left\{ -\frac{\sum_{n=1}^N (x_n - \mu)^2}{2\sigma^2} \right\}$$

$$\Rightarrow \text{Posterior distribution } P(\mu | X) \propto P(X | \mu) \cdot P(\mu) \\ = N(\mu | \mu_N, \sigma_N^2)$$

$$\mu_N = \frac{\sigma^2}{N\sigma_0^2 + \sigma^2} \mu_0 + \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2} \cdot \mu_{ML}$$

$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}$$

Case 2: known mean, unknown variance.

$$\Rightarrow \text{likelihood: } P(X | \mu) = \prod_{n=1}^N N(x_n | \mu, \lambda^{-1})$$

$$\propto \lambda^{\frac{N}{2}} \exp \left\{ -\frac{\lambda}{2} \sum_{n=1}^N (x_n - \mu)^2 \right\}$$

then Conjugate prior: In case of using $\text{Gam}(\lambda | a_0, b_0)$
 we have $P(\lambda | x) \propto \lambda^{a_0+1} \lambda^{\frac{N}{2}} \exp \left\{ -b_0 \lambda - \frac{\lambda}{2} \sum_{n=1}^N (x_n - \mu)^2 \right\}$
 $= \text{Gam}(\lambda | a_N, b_N)$

where $a_N = a_0 + \frac{N}{2}$, $b_N = b_0 + \frac{N}{2} S_{\mu}^2$

Case 3: unknown mean, unknown variance.

\Rightarrow likelihood: $P(X | \mu, \lambda) = \prod_{n=1}^N \left(\frac{\lambda}{2\pi} \right)^{\frac{1}{2}} \exp \left\{ -\frac{\lambda}{2} (x_n - \mu)^2 \right\}$
 $\propto \left[\lambda^{\frac{1}{2}} \exp \left(-\frac{\lambda \mu^2}{2} \right) \right]^N \exp \left\{ \lambda \mu \cdot \sum_{n=1}^N x_n - \frac{\lambda}{2} \sum_{n=1}^N x_n^2 \right\}$

Conjugate prior: $P(\mu, \lambda) \propto \left[\lambda^{\frac{1}{2}} \exp \left(-\frac{\lambda \mu^2}{2} \right) \right]^{\beta}$
 $\exp [c \lambda \mu - d \lambda]$
 $= \exp \left\{ -\frac{\beta \lambda}{2} \left(\mu - \frac{c}{\beta} \right)^2 \right\} \lambda^{\beta/2}$

~~Student's~~ Student's t-distribution \leftarrow normal $\exp \left\{ -\left(d - \frac{c^2}{2\beta} \right) \lambda \right\}$

$$P(X | \mu, a, b) = \int_0^{\infty} N(X | \mu, \tau^{-1}) \text{Gam}(\tau | a, b) d\tau$$

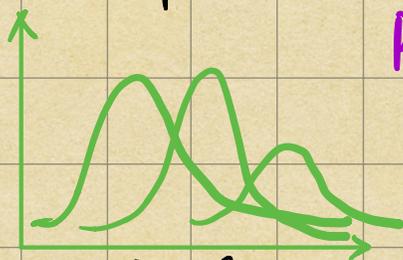
$$= \text{St}(X | \mu, \lambda, \nu) \text{ where } \lambda \text{ is precision \&}$$

when $\nu \rightarrow \infty$:

ν is degree of freedom.

$$\text{St}(X | \mu, \lambda, \nu) \rightarrow N(X | \mu, \lambda^{-1})$$

★ mixture of Gaussians:



$$P(X) = \sum_{k=1}^K \pi_k N(X | \mu_k, \Sigma_k), \text{ where } \sum_{k=1}^K \pi_k = 1$$

$$= \sum_{k=1}^K P(k) \cdot P(X | k) \quad \& \pi_k \in [0, 1]$$

Let $\Lambda = \{\pi\} = \{\pi_1, \dots, \pi_k\}$, $\mu = \{\mu_1, \dots, \mu_k\}$, $\Sigma = \{\Sigma_1, \dots, \Sigma_k\}$

★ The Exponential Family (brood class of disturb)

$$P(X | \eta) = h(x) g(\eta) \exp \{ \eta^T u(x) \}$$

13. Example 1: Bernoulli distrib

$$P(x|\mu) = \mu^x (1-\mu)^{1-x} = \exp\{x \ln \mu + (1-x) \cdot \ln(1-\mu)\}$$

$$= (1-\mu) \exp\left\{\ln \frac{\mu}{1-\mu} \cdot x\right\}$$

$\mu = \frac{1}{1 + \exp(-\eta)} = \sigma(\eta)$
 $\downarrow \eta$
 $\downarrow u(x)$

Example 2: Gaussian distrib

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

$$= \frac{1}{(\sqrt{2\pi}\sigma)^2} \exp\left\{-\frac{1}{2\sigma^2}x^2 + \frac{\mu}{\sigma^2}x - \frac{1}{2\sigma^2}\mu^2\right\}$$

$$= h(x) \cdot g(\eta) \cdot \exp\{\eta^T u(x)\}$$

where $\eta = \begin{pmatrix} \frac{\mu}{\sigma^2} \\ -\frac{1}{2\sigma^2} \end{pmatrix}$ $u(x) = \begin{pmatrix} x \\ x^2 \end{pmatrix}$ $h(x) = (\sqrt{2\pi})^{-2}$ $g(\eta) = (\sqrt{2\pi})^2 \exp\left(\frac{\eta^T \eta}{2}\right)$

Conjugate Prior:

$P(\eta|x, v) = t(x, v) \cdot g(\eta)^v \exp\{v \eta^T x\}$

↻ Conjugate Prior normalization coeff

$P(\eta|x, \mathcal{X}, v) \propto g(\eta)^{v+N} \exp\left\{\eta^T \left(\sum_{n=1}^N u(x_n) + v \mathcal{X}\right)\right\}$

↻ Posterior

Day 3: Linear Models for regression.

3.1 Linear Basis Function Models: Basis function $\{\phi_j(x)\} \equiv$ feature

$P(t|x)$ is the uncertainty of t for certain x .

$$y(x, w) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(x) = \sum_{j=0}^{M-1} w_j \phi_j(x) = w^T \phi(x)$$

\downarrow bias parameter
 \downarrow basis function
 $w = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_{M-1} \end{bmatrix}$

$\phi_j(x) = x^j$ } kernel function $k(x_n, x_m) = \phi^T(x_n) \cdot \phi(x_m)$

4. ~~①~~ $\phi_j(x) = \exp\left\{-\frac{(x-\mu_j)^2}{2\sigma^2}\right\}$ (实验效果较差)

↑ Gaussian basis function with spatial.

② sigmoid basis function

$$\phi_j(x) = \delta\left(\frac{x-\mu_j}{s}\right) \& \delta(a) = \frac{1}{1+e^{-a}}$$

$$\tanh(a) = 2\delta(a) - 1 = \frac{1-e^{-a}}{1+e^{-a}}$$

Maximum likelihood & least squares

$$t = y(x, w) + \epsilon \quad \epsilon \sim N(0, \beta^{-1}) \rightarrow \text{Gaussian Noise}$$

$$\rightarrow P(t|x, w, \beta) = N(t|y(x, w), \beta^{-1})$$

$$\rightarrow E[t|x] = \int t p(t|x) \cdot dt = y(x, w)$$

Input data $X = \{x_1, x_2, \dots, x_N\}$

$$P(t|x, w, \beta) = \prod_{n=1}^N N(t_n | w^T \phi(x_n), \beta^{-1})$$

$$\ln P(t|x, w, \beta) = \sum_{n=1}^N \ln N(t_n | w^T \phi(x_n), \beta^{-1})$$

$$= \sum_{n=1}^N \ln \beta - \frac{N}{2} \ln(2\pi) - \beta \left\{ \frac{1}{2} \sum_{n=1}^N |t_n - w^T \phi(x_n)|^2 \right\}$$

$$\nabla_w \ln P(t|x, w, \beta) = \sum_{n=1}^N \{t_n - w^T \phi(x_n)\} \phi^T(x_n) \stackrel{\text{let}}{=} 0$$

$$\therefore \sum_{n=1}^N t_n \phi(x_n)^T - w^T \left(\sum_{n=1}^N \phi(x_n) \cdot \phi(x_n)^T \right) = 0$$

$$w_{ML} = (\Phi^T \Phi)^{-1} \Phi^T t, \text{ where } \Phi = \begin{pmatrix} \phi_0(x_1) & \phi_1(x_1) & \dots & \phi_{M-1}(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & \dots & \phi_{M-1}(x_2) \\ \vdots & \vdots & \ddots & \vdots \end{pmatrix}$$

The error function becomes

$$E_D(w) = \frac{1}{2} \sum_{n=1}^N \left\{ t_n - w_0 - \sum_{j=1}^{M-1} w_j \phi_j(x_n) \right\}^2$$

$$\text{let } \frac{\partial E_D(w_0)}{\partial w_0} = 0, \therefore w_0 = \frac{1}{N} \sum_{n=1}^N t_n - \sum_{j=1}^{M-1} w_j \left(\frac{1}{N} \sum_{n=1}^N \phi_j(x_n) \right)$$

Error function could be seen as a maximum likelihood solution with Gaussian Noise Model.

15.

Sequential Learning :

★ Stochastic Gradient Descent (SGD)

$$W^{(L+1)} = W^{(L)} - \eta \nabla E_n$$

$$= W^{(L)} - \eta (t_n - W^{(L)\top} \phi_n) \phi_n$$

Regularization Least Square :

$$E_w(w) = \frac{1}{2} w^\top w = \frac{1}{2} \|w\|^2$$

(weight decay) η (hyperparameter)

$$w_{RLS} = (\lambda I + \Phi^\top \Phi)^{-1} \Phi^\top t$$

3.2 The Bias-Variance Decomposition.

The expected squared loss can be written by:

$$E[L] = \int \{y(x) - h(x)\}^2 p(x) dx + \int \{h(x) - t\}^2 \cdot P(x,t) dx dt$$

expected loss = (bias)² + variance + noise.

average prediction : $\bar{y}(x) = \frac{1}{N} \sum_{l=1}^N y^{(l)}(x)$ prediction function

$$(bias)^2 = \frac{1}{N} \sum_{n=1}^N \{ \bar{y}(x_n) - h(x_n) \}^2$$

$$variance = \frac{1}{N} \sum_{n=1}^N \frac{1}{N} \sum_{l=1}^N \{ y^{(l)}(x_n) - \bar{y}(x_n) \}^2$$

3.3 Bayesian Linear Regression

Given w conjugate prior

$$P(w) = N(w | m_0, S_0)$$

$$P(w|t) \propto P(t|w) \cdot P(w) = N(w | m_N, S_N)$$

↑ prior

$$\Rightarrow m_N = S_N (S_0^{-1} m_0 + \beta \Phi^\top t) = w_{MAP}$$

To simplify the treatment.

$$S_N^{-1} = S_0^{-1} + \beta \Phi^\top \Phi$$

we consider a zero-mean isotropic Gaussian as a prior

$$P(w|\alpha) = N(w | 0, \alpha^{-1} I)$$

$$\Rightarrow m_N = \beta S_N \Phi^\top t \quad \& \quad \ln P(w|t) = -\frac{\beta}{2} \sum_{n=1}^N \{ t_n - w^\top \phi(x_n) \}^2 - \frac{\alpha}{2} w^\top w$$

$$S_N^{-1} = \alpha I + \beta \Phi^\top \Phi$$

$$\therefore P(W|\alpha) = \left[\frac{\Gamma(\frac{\alpha}{2})}{\Gamma(\frac{1}{2})} \right]^\alpha \exp\left(-\frac{\alpha}{2} \sum_{j=1}^N |w_j|^2\right)$$

★ Prediction distribution

$$P(t|t, \alpha, \beta) = \int P(t|w, t, \beta) P(w|t, \alpha, \beta) dw$$

↑ Gaussian $N(t|y(x, w), \beta^{-1})$

$$\star P(t|x, t, \alpha, \beta) = N(t | m_N^T \phi(x), \sigma_N^2(x))$$

$$\text{where } \sigma_N^2(x) = \frac{1}{\beta} + \phi(x)^T S_N \phi(x)$$

noise data ↓ uncertainty due to w

Equivalent Kernel:

$$\star y(x, m_N) = m_N^T \phi(x) = \beta \phi^T(x) \cdot S_N \Phi^T(t)$$

$$= \sum_{n=1}^N \beta \phi^T(x) \cdot S_N \phi(x_n, t_n) = \sum_{n=1}^N k(x, x_n) \cdot t_n$$

where $k(x, x') = \beta \cdot \phi^T(x) S_N \phi(x')$ ○ prediction

★ Bayesian Model Comparison: set a criterion to select model.

$$P(M_i | D) \propto P(M_i) \cdot P(D | M_i)$$

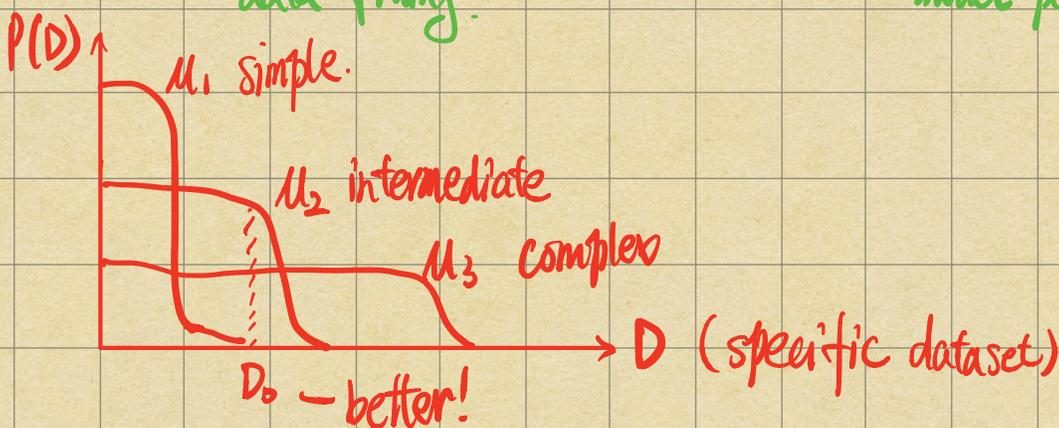
$$P(D | M_i) = \int P(D | w, M_i) P(w | M_i) dw$$

$$\star P(D) = \int P(D | w) \cdot p(w) \cdot dw \approx P(D | w_{MAP}) \cdot \frac{\Delta W_{\text{posterior}}}{\Delta W_{\text{prior}}}$$

↓ ΔW_{prior} ↓ $\Delta W_{\text{posterior}}$ assume sharp.

$$\Rightarrow \ln P(D) \approx \ln P(D | w_{MAP}) + \ln \left(\frac{\Delta W_{\text{posterior}}}{\Delta W_{\text{prior}}} \right)$$

data fitting ← model penalty.



17. 3.5 The Evidence Approximation.

cross validation: for deciding hyperparameter.

Predictive distribution:

$$P(t|t) = \int \int P(t|\omega, \beta) \cdot P(\omega|t, \alpha, \beta) \cdot P(\alpha, \beta|t) d\omega d\alpha d\beta$$

$$\star P(t|t) \approx P(t|t, \hat{\alpha}, \hat{\beta}) = \int P(t|\omega, \hat{\beta}) \cdot P(\omega|t, \hat{\alpha}, \hat{\beta}) \cdot d\omega$$

$$(\hat{\alpha}, \hat{\beta}) = \underset{(\alpha, \beta)}{\operatorname{argmax}} P(\alpha, \beta|t) = \underset{(\alpha, \beta)}{\operatorname{argmax}} P(t|\alpha, \beta) \cdot P(\alpha, \beta)$$

~~★~~ Evaluation of evidence function \rightarrow predictive distribution.

$$P(t|\alpha, \beta) = \int P(t|\omega, \beta) \cdot P(\omega|\alpha) \cdot d\omega$$

$$= \left(\frac{\beta}{2\pi}\right)^{N/2} \left(\frac{\alpha}{2\pi}\right)^{N/2} \int \exp\{-E(\omega)\} \cdot d\omega.$$

where $E(\omega) = \beta E_p(\omega) + \alpha E_u(\omega)$

$$= \frac{\beta}{2} \|t - \Phi \omega\|^2 + \frac{\alpha}{2} \omega^T \omega.$$

★ \int is for evidence!!!

$$= \left[E(m_N) \right] + \left[\frac{1}{2} (\omega - m_N)^T A (\omega - m_N) \right]$$

$$\left[\frac{\beta}{2} \|t - \Phi m_N\|^2 + \frac{\alpha}{2} m_N^T m_N \right] \left[\alpha I + \beta \Phi^T \Phi \right]$$

$\downarrow E_p(\omega)$ $\downarrow E_u(\omega)$

$$\int \exp\{-E(\omega)\} d\omega = \exp\{-E(m_N)\} \int \exp\left\{-\frac{1}{2} (\omega - m_N)^T A (\omega - m_N)\right\} \cdot d\omega$$

$$= \exp\{-E(m_N)\} (2\pi)^{N/2} |A|^{-\frac{1}{2}}$$

$$\ln P(t|\alpha, \beta) = \frac{N}{2} \ln \beta + \frac{N}{2} \ln \alpha - \frac{N}{2} \ln(2\pi) - E\{m_N\} - \frac{1}{2} \ln |A|$$

Maximizing the evidence function

$$\textcircled{1} \hat{\alpha} = \underset{\alpha}{\operatorname{argmax}} P(t|\alpha, \beta)$$

Given the eigenvector equation $(\beta \Phi^T \Phi) u_i = \lambda_i u_i$
 A has eigenvalues $\{\alpha + \lambda_i\}$

18.

$$\frac{d}{d\alpha} \ln |A| = \frac{d}{d\alpha} \ln \prod_i (\lambda_i + \alpha) = \frac{d}{d\alpha} \sum_i \ln(\lambda_i + \alpha) = \sum_i \frac{1}{\lambda_i + \alpha}$$

$$\therefore \frac{d}{d\alpha} \ln P(t|\alpha, \beta) = \frac{M}{2\alpha} - \frac{1}{2} m_N^T \cdot m_N - \frac{1}{2} \sum_i \frac{1}{\lambda_i + \alpha} = 0$$

$$\alpha m_N^T \cdot m_N = M - \alpha \sum_{i=1}^M \frac{1}{\lambda_i + \alpha} = r = \sum \frac{\lambda_i}{\alpha + \lambda_i}$$

$$\star \alpha = \frac{r}{m_N^T \cdot m_N} \quad \star \alpha^{(0)} \rightarrow m_N^{(0)} \rightarrow r^{(0)} \rightarrow \alpha^{(1)} \rightarrow \dots$$

$$\textcircled{2} \hat{\beta} = \arg \max P(t|\alpha, \beta)$$

$$\frac{d}{d\beta} \ln |A| = \frac{d}{d\beta} \sum_i \ln(\lambda_i + \alpha) = \frac{1}{\beta} \sum_i \frac{\lambda_i}{\lambda_i + \alpha} = \frac{r}{\beta}$$

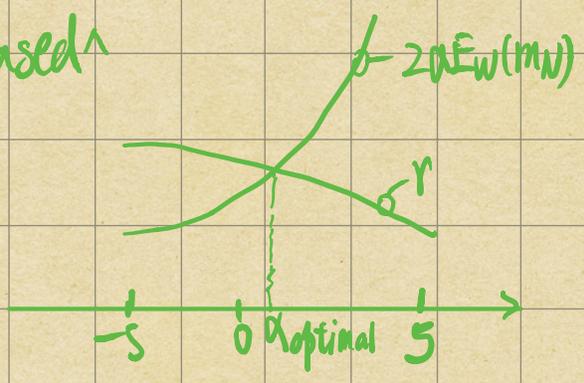
$$\frac{d}{d\beta} \ln P(t|\alpha, \beta) = \frac{N}{2\beta} - \frac{1}{2} \sum_{n=1}^N \{t_n - m_N^T \phi(x_n)\}^2 - \frac{r}{2\beta} = 0$$

$$\star \beta^{-1} = \frac{1}{N-r} \sum_{n=1}^N \{t_n - m_N^T \phi(x_n)\}^2 \quad \star \beta^{(0)} \rightarrow m_N^{(0)} \rightarrow r^{(0)} \rightarrow \beta^{(1)} \rightarrow \dots$$

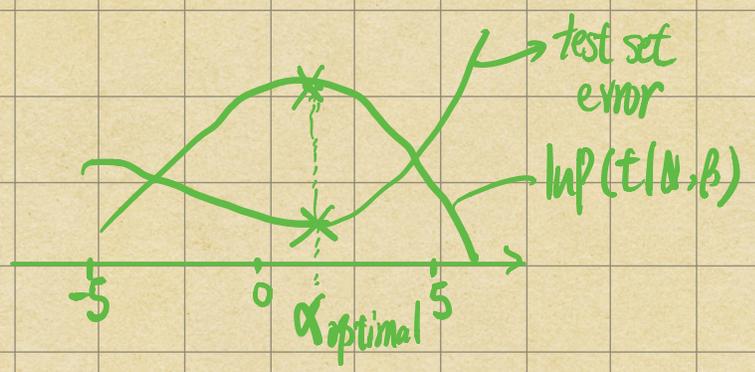
effective number of parameters

$$0 \leq \frac{\lambda_i}{\lambda_i + \alpha} \leq 1 \quad 0 \leq r \leq M$$

biased



unbiased



activate function

19. Day 4: Linear Models for Classification $y(x) = f(w^T x + w_0)$

intro: 1. Generative Model $\rightarrow P(x|C_i)$

* Discriminative Model $\rightarrow P(C_i|x)$

3. Linear Discriminant Function

4. Discriminant Function:

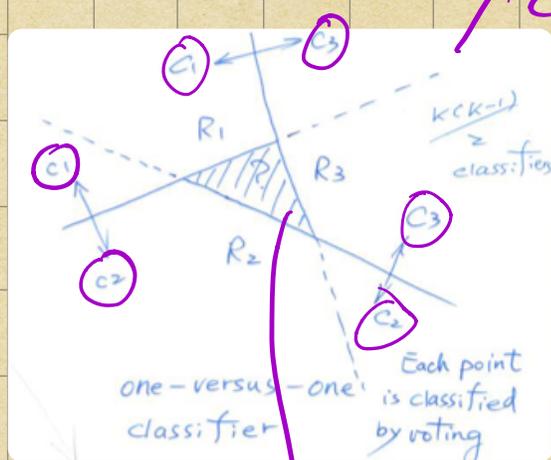
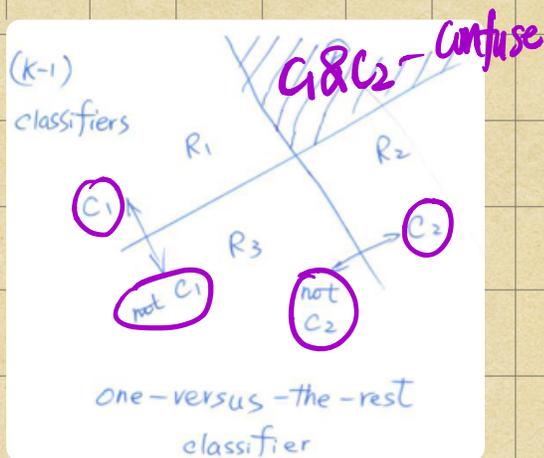
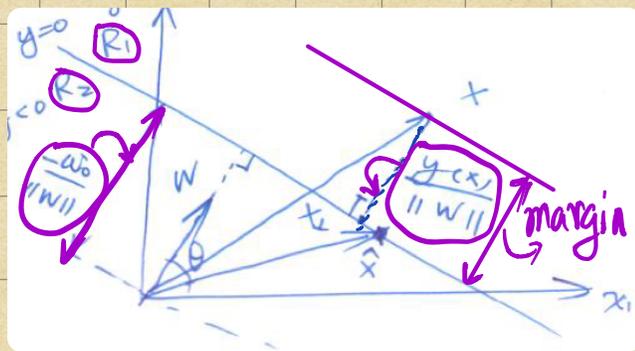
two classes: $y(x) = w^T x + w_0$

$\frac{w}{\|w\|} = x_{\perp}$ $x = x_{\perp} + r \frac{w}{\|w\|}$
 $\Rightarrow w^T x = w^T x_{\perp} + r \frac{w^T w}{\|w\|}$

with $w^T x_{\perp} + w_0 = 0 \quad \therefore w^T x = -w_0 + r \frac{w^T w}{\|w\|}$

$\Rightarrow r = \frac{y(x) - w_0}{\|w\|}$ *

Multiple classes ($K > 2$)

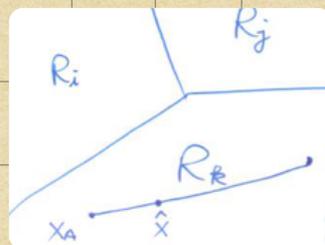


we can avoid these difficulties by:

$y_k(x) = w_k^T x + w_{k0}$

Decision boundary between C_k & C_j is given by

$y_k(x) = y_j(x) \Rightarrow (w_k - w_j)^T x + (w_{k0} - w_{j0}) = 0$



① Least squares for classification

$$y_k(x) = w_k^T x + w_{k0} \quad k=1, 2, \dots, K \text{ (class) [linear discriminant function]}$$

Error function can be written by:

$$E_D(\tilde{w}) = \frac{1}{2} \text{Tr} \{ (\tilde{X} \tilde{w} - T)^T (\tilde{X} \tilde{w} - T) \}$$

$$\text{let } \nabla_w E_D(\tilde{w}) = 0 \Rightarrow \tilde{w} = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T T = \tilde{X}^+ T$$

$$y(x) = \tilde{w}^T \tilde{x} = T^T (\tilde{X}^+)^T \tilde{x}$$

② Linear Discriminant (aka: LDA \leftrightarrow PCA Linear Discriminant Analysis)

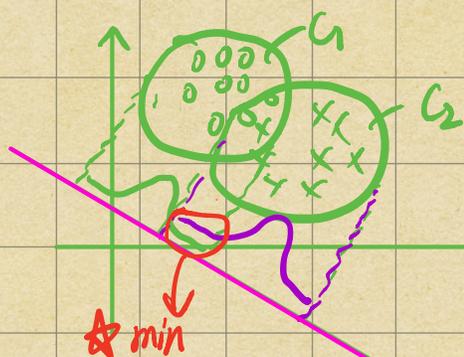
Dimensionality Reduction: \rightarrow projection.

$$y = w^T x \begin{cases} \geq c_1 \\ < c_2 \end{cases} w_0$$

Fisher's Criterion: $J(w) = \frac{\text{Mean} \downarrow (m_2 - m_1)^2}{\text{var} \rightarrow S_1^2 + S_2^2} = \frac{w^T S_B w}{w^T S_W w}$

$$m_1 = \frac{1}{N} \sum_{n \in C_1} x_n, \quad m_2 = \frac{1}{N_2} \sum_{n \in C_2} x_n \quad m_2 - m_1 = w^T (m_2 - m_1)$$

$$\text{where } S_B = (m_2 - m_1)(m_2 - m_1)^T \quad S_W = \sum_{n \in C_1} (x_n - m_1)(x_n - m_1)^T + \sum_{n \in C_2} (x_n - m_2)(x_n - m_2)^T$$



- ① less variance
- ② more distance

$$\text{let } \nabla_w J(w) = 0 \Rightarrow (w^T S_B w) S_W w = (w^T S_W w) S_B w$$

$$\text{where } S_B w = (m_2 - m_1)(m_2 - m_1)^T \cdot w$$

$$\text{finally, } w \propto S_W^{-1} (m_2 - m_1)$$

③ perceptron algorithm

21.

4.2 Probabilistic Generative Models

→ divide $P(X|C_1)P(C_1)$

$$P(C_1|X) = \frac{P(X|C_1) \cdot P(C_1)}{P(X|C_1) \cdot P(C_1) + P(X|C_2) \cdot P(C_2)} = \frac{1}{1 + \exp(-a)} = \sigma(a)$$

where $a = \ln \frac{P(X|C_1) P(C_1)}{P(X|C_2) P(C_2)}$

logistic sigmoid function

For $k > 2$ $P(C_k|X) = \frac{P(X|C_k) \cdot P(C_k)}{\sum_j P(X|C_j) \cdot P(C_j)} = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$

Generative Model: $P(X|C_i)$ [Indirect Model]

where $a_k = \ln P(X|C_k) P(C_k)$

Continuous Inputs

$$P(X|C_k) = \frac{1}{(2\pi)^{D/2}} \cdot \frac{1}{|\Sigma|^{D/2}} \exp \left\{ -\frac{1}{2} (X - \mu_k)^T \Sigma^{-1} (X - \mu_k) \right\}$$

$$P(C_i|X) = \sigma(w^T X + w_0) \quad \{ \mu_1, \mu_2, \Sigma, X = P(C_j) \}$$

$$w = \Sigma^{-1} (\mu_1 - \mu_2) \quad w_0 = -\frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2 + \ln \frac{P(C_1)}{P(C_2)}$$

Maximum likelihood solution:

Gaussian class-conditional density.

$$P(C_1) = \pi \quad P(C_2) = 1 - \pi \quad t_n = 0/1$$

$$P(X_n, C_1) = P(C_1) \cdot P(X_n|C_1) = \pi N(X_n|\mu_1, \Sigma) \quad t_n = 1$$

$$P(X_n, C_2) = P(C_2) \cdot P(X_n|C_2) = (1 - \pi) \cdot N(X_n|\mu_2, \Sigma) \quad t_n = 0.$$

likelihood function:

$$P(t_1 \pi, \mu_1, \mu_2, \Sigma) = \prod_{n=1}^N [\pi N(X_n|\mu_1, \Sigma)]^{t_n} [(1 - \pi) N(X_n|\mu_2, \Sigma)]^{1-t_n}$$

After lots of compute:

$$1. \pi = \frac{1}{N} \sum_{n=1}^N t_n = \frac{N_1}{N} = \frac{N_1}{N_1 + N_2}$$

$$2. \mu_1 = \frac{1}{N_1} \sum_{n=1}^N t_n \cdot X_n$$

$$3. \mu_2 = \frac{1}{N_2} \sum_{n=1}^N (1 - t_n) \cdot X_n$$

$$4. \Sigma = -\frac{N}{2} \ln |\Sigma| - \frac{N}{2} \text{Tr} \{ \Sigma^{-1} s \}$$

22.

$$S_1 = \frac{1}{N_1} \sum_{n \in C_1} (x_n - \mu_1)(x_n - \mu_1)^T \Rightarrow S = \frac{N_1}{N} S_1 + \frac{N_2}{N} S_2$$

$$S_2 = \frac{1}{N_2} \sum_{n \in C_2} (x_n - \mu_2)(x_n - \mu_2)^T \quad \Sigma = S$$

★★★
★★★
必考!

4.3 Probabilistic Discriminative Models (Direct Model) $P(C_k|X)$

Intro: Indirect: $P(X|C_k) \rightarrow ML \rightarrow P(C_k|X) \rightarrow P(C_k)$

Direct: $P(C_k|X) \rightarrow$

$$P(C_1|\phi) = y(\phi) = \sigma(w^T \phi) \quad \frac{dy}{da} = \sigma(1-y)$$

$$P(C_2|\phi) = 1 - P(C_1|\phi)$$

Cross-Entropy

$$D = \{ \phi_n, t_n \}_{n=1}^N, \quad t_n \in \{0, 1\}$$

$$E(w) = -\ln P(t|w)$$

$$P(t|w) = \prod_{n=1}^N y_n^{t_n} (1-y_n)^{1-t_n}$$

$$= -\sum_{n=1}^N \{ t_n \ln y_n + (1-t_n) \ln (1-y_n) \}$$

Iterative Reweighted Least Squares (IRLS)

1. $w^{(new)} = w^{(old)} - \eta \nabla E(w)$ steepest descent alg.

2. Newton-Raphson alg

$$w^{(new)} = w^{(old)} - H^{-1} \nabla E(w), \quad \text{where } H = \nabla \nabla E(w)$$

$$\nabla E(w) = \sum_{n=1}^N (y_n - t_n) \phi_n = \Phi^T (y - t) \quad \text{data matrix}$$

$$H = \nabla \nabla E(w) = \sum_{n=1}^N y_n (1-y_n) \phi_n \phi_n^T = \Phi^T R \Phi, \quad \text{weighting matrix } R_{NN} = [R_{nn}]$$

$$\Rightarrow w^{(new)} = w^{(old)} - (\Phi^T R \Phi)^{-1} \Phi^T (y - t)$$

$$= (\Phi^T R \Phi)^{-1} \{ \Phi^T R \Phi w^{(old)} - \Phi^T (y - t) \}$$

$$= (\Phi^T R \Phi)^{-1} \Phi^T R z, \quad \text{where } z = \Phi w^{(old)} - R^{-1} (y - t)$$

core idea: $w \rightarrow R \rightarrow w \rightarrow R$

4.4 Laplace Approximation

classification $\rightarrow \int S \rightarrow$ 不可积 \Rightarrow approximation. (Taylor & ln)

$$P(z) = \frac{f(z)}{\int f(z) dz} \quad \text{via Taylor series to quadratic equation.}$$

$$\ln f(z) \approx \ln f(z_0) - \frac{1}{2} (z - z_0)^T A (z - z_0), \quad \text{where } A = -\nabla \nabla \ln f(z)|_{z=z_0}$$

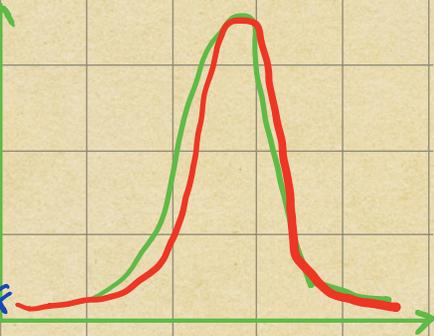
$$\Rightarrow f(z) \approx f(z_0) \exp \left\{ -\frac{1}{2} (z - z_0)^T A (z - z_0) \right\}$$

$$\therefore \int f(z) dz \approx f(z_0) \frac{(2\pi)^{M/2}}{|A|^{1/2}} \quad * \quad \rightarrow$$

$$P(z) \leftarrow q(z)$$

↑ non-Gaussian.
↑ Gaussian

$$q(z) = \frac{|A|^{1/2}}{(2\pi)^{M/2}} \exp \left\{ -\frac{1}{2} (z - z_0)^T A (z - z_0) \right\} * *$$



Model comparison and BIC,

We have model evidence: $P(D) = \int P(D|\theta) \cdot P(\theta) d\theta$

$$\frac{P(D|\theta) \cdot P(\theta)}{P(D)} = P(\theta|D) = \frac{f(\theta)}{\int f(\theta) d\theta}$$

$$\ln P(D) \approx \ln P(D|Q_{MAP}) + \ln P(Q_{MAP}) + \frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln |A|$$

Day 5: kernel function: Memory based Methods

$$K(x, x') = \phi(x)^T \cdot \phi(x') \quad * \quad \text{similarity between } x \text{ \& } x'$$

Gaussian Process (可能考)

6.1: Dual Representation

$$J(w) = \frac{1}{2} \sum_{n=1}^N \{ w^T \phi(x_n) - t_n \}^2 + \frac{\lambda}{2} w^T w$$

$$\nabla_w J(w) = 0 \rightarrow w_{opt} = -\frac{1}{\lambda} \sum_{n=1}^N \{ w^T \phi(x_n) - t_n \} \phi(x_n) = \sum_{n=1}^N a_n \phi(x_n) = \Phi^T a$$

$$a_n = -\frac{1}{\lambda} \{ w^T \phi(x_n) - t_n \}, \text{ where } \Phi^T = [\phi(x_1), \dots, \phi(x_N)], a = [a_1, \dots, a_N]^T$$

$$\Rightarrow J(a) = \frac{1}{2} a^T \Phi \Phi^T \Phi \Phi^T a - a^T \Phi \Phi^T t + \frac{1}{2} t^T t + \frac{\lambda}{2} a^T \Phi \Phi^T a$$

let $K = \Phi \cdot \Phi^T$, we obtain $K_{nm} = \phi(x_n)^T \phi(x_m) = K(x_n, x_m)$

$$J(a) = \frac{1}{2} a^T K K a - a^T K t + \frac{1}{2} t^T t + \frac{\lambda}{2} a^T K a$$

$$\nabla_a J(a) = 0 \quad a_{opt} = (K + \lambda I_N)^{-1} t$$

$$y(x) = w^T \phi(x) = a^T \Phi \phi(x) = K(x)^T (K + \lambda I_N)^{-1} t$$

6.4: Gaussian Process

Linear regression revisited:

$$y(x) = w^T \phi(x) \quad \& \quad p(w) = N(w | 0, \alpha^{-1} I)$$

Given $\{x_1, \dots, x_n\}$ we have $y = \{y(x_1), \dots, y(x_n)\}^T$

$$\Rightarrow y = \Phi w. \quad \therefore E[y] = \Phi E[w] = 0$$

$$\text{cov}[y] = E[yy^T] = \Phi E[w w^T] \Phi^T = \frac{1}{\alpha} \Phi \Phi^T = K.$$

★ ★ Gaussian Process for regression:

$$t_n = y_n + \epsilon_n \text{ noise}$$

$$P(t_n | y_n) = N(t_n | y_n, \beta^{-1})$$

$$P(t | y) = N(t | y, \beta^{-1} I_N) \text{ or } N(t - y | 0, \beta^{-1} I_N)$$

$$P(y) = N(y | 0, k)$$

$$P(t) = \int P(t | y) \cdot P(y) dy = N(t | 0, C_{t,t})$$

$$\text{where } C(x_n, x_m) = k(x_n, x_m) + \beta^{-1} \delta_{nm}$$

One widely used kernel function:

$$k(x_n, x_m) = \theta_0 \exp\left\{-\frac{\theta_1}{2} \|x_n - x_m\|^2\right\} + \theta_2 + \theta_3 x_n^T x_m$$

$(\theta_1, \theta_2, \theta_3, \theta_4)$ & $P(t)$ hyperparameter.

$$P(t_{n+1}) = P(t_1, \dots, t_{n+1}) = N(t_{n+1} | 0, C_{n+1}), \text{ where } C_{n+1} = \begin{pmatrix} C_N & k \\ k^T & C \end{pmatrix}$$

$$C = k(x_{n+1}, x_{n+1}) + \beta^{-1}$$

$$P(t_{n+1} | t) = N(t_{n+1} | m(x_{n+1}), \delta^2(x_{n+1})) = K^T C_N^{-1} t = C - K^T C_N^{-1} K$$

Learning parameters:

$$\hat{\theta} = \text{argmax}_{\theta} \log P(t | \theta) \rightarrow \ln P(t | \theta) = -\frac{1}{2} \ln |C_N| - \frac{1}{2} t^T C_N^{-1} t - \frac{N}{2} \ln(2\pi)$$

$$\frac{\partial}{\partial \theta_i} \ln P(t | \theta) = -\frac{1}{2} \text{Tr}(C_N^{-1} \frac{\partial C_N}{\partial \theta_i}) + \frac{1}{2} t^T C_N^{-1} \frac{\partial C_N}{\partial \theta_i} C_N^{-1} t$$

25. 不考. 6.5 Gaussian processes for classification

太复杂

$$P(t|a) = \delta(a)^T (1 - \delta(a))^T$$

Laplace approximation:

$$P(a_{n+1}|t_n) = \int P(a_{n+1}, a_n | t_n) da_n P(a_{n+1}|a_n) P(a_n|t_n) \cdot P(t_n|a_n)$$

$$= \int P(a_{n+1}|a_n) P(a_n|t_n) da_n$$

$$q(a_n) = N(a_n | a_n^*, H^{-1})$$

$$\rightarrow P(a_{n+1}|a_n) \cong \int P(a_{n+1}|a_n) q(a_n) da_n$$

$$\text{then } P(t_{n+1}=1|t_n) = \int P(t_{n+1}=1|a_{n+1}) P(a_{n+1}|t_n) da_{n+1}$$

Day 6. 7: Sparse Kernel Machines (linear model)

Intro: 1. Support Vector Machines (SVM) \Rightarrow (select important vector to predict)
 2. Relevance Vector Machines (RVM)

Support Vector Machines.
 Classification:

(-): Nonoverlapping classes

(=): Overlapping classes

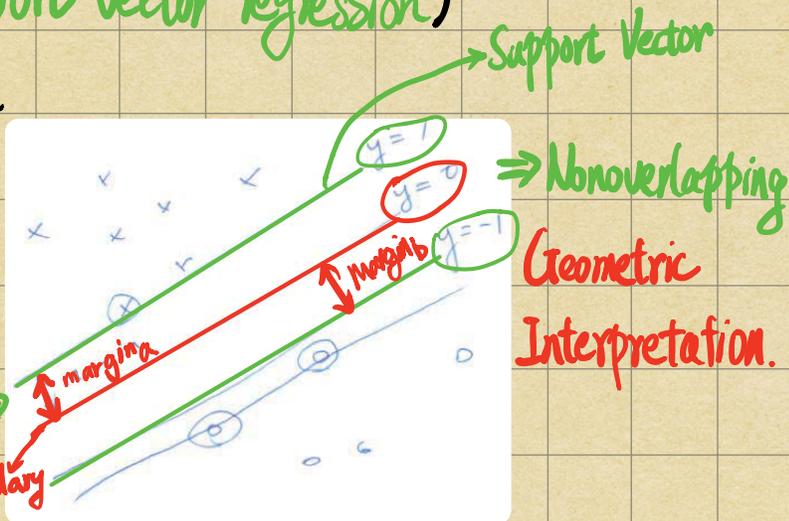
Regression: (support vector regression)

Two-class classification problem

$$y = w^T \phi(x) + b$$

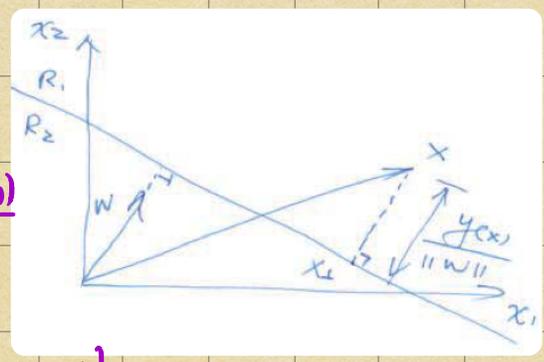
found a bound maximum the "Margin" among dataset.

Margin: $\text{margin}_a = \text{margin}_b$



26.

Distance: $v = \frac{y(x)}{\|w\|}$



add class label we obtain:

$\Rightarrow \text{margin} = \frac{t_n y(x_n)}{\|w\|} = \frac{t_n (w^T \phi(x_n) + b)}{\|w\|}$

★ Maximum margin solution is:

$(\hat{w}, \hat{b}) = \underset{(w, b)}{\text{argmax}} \left\{ \frac{1}{\|w\|} \min_n [t_n (w^T \phi(x_n) + b)] \right\}$

with $t_n (w^T \phi(x_n) + b) \geq 1$ (as constrain) we obtain: $(\hat{w}, \hat{b}) = \underset{(w, b)}{\text{argmin}} \frac{1}{2} \|w\|^2$

$\Rightarrow L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{n=1}^N \alpha_n [t_n (w^T \phi(x_n) + b) - 1]$

let $\frac{\partial L}{\partial w} = 0, \frac{\partial L}{\partial b} = 0$ Lagrange multiplier

$\therefore w = \sum_{n=1}^N \alpha_n t_n \phi(x_n) \quad 0 = \sum_{n=1}^N \alpha_n t_n$

Dual representation of the maximum margin problem.

$\rightarrow \max_{\alpha} \left\{ \tilde{L}(\alpha) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m t_n t_m K(x_n, x_m) \right\}$, where $\begin{cases} \sum_{n=1}^N \alpha_n t_n = 0 \\ \alpha_n \geq 0 \\ t_n y(x_n - 1) \geq 0 \end{cases}$
Solution: Sequential minimal optimization (SMO)

In classification problem:

$y(x) = \sum_{n=1}^N \alpha_n t_n k(x, x_n) + b$
↑ test data ↑ training data

the KKT conditions: $\alpha_n \geq 0, t_n y(x_n) - 1 = 0 \quad \alpha_n [t_n y(x_n) - 1] = 0$

$\therefore \alpha_n = 0$ or $t_n y(x_n) = 1$

the parameter b:

$t_n \left(\sum_{m \in S} \alpha_m t_m k(x_n, x_m) + b \right) = 1$
 $\Rightarrow b = \frac{1}{N_S} \sum_{n \in S} (t_n - \sum_{m \in S} \alpha_m t_m k(x_n, x_m))$

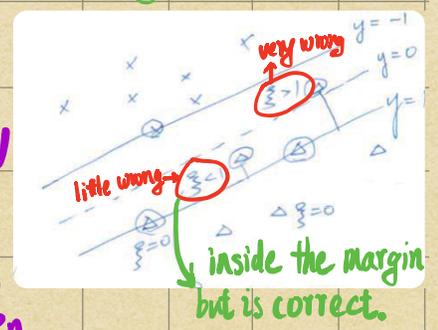
error function: $\sum_{n=1}^N E_n (y(x_n) t_n - 1) + \lambda \|w\|^2$ → quadratic regularizer.

Overlapping class distribution:

Slack variables are introduced: $\xi_n \geq 0, n=1, \dots, N$

↓ to measure to misclassified point.

classification constraints are replaced by: $t_n y(x_n) \geq 1 - \xi_n$



27.

Therefore, $\min_w \left(C \sum_{n=1}^N \xi_n + \frac{1}{2} \|w\|^2 \right)$
 \downarrow hyperparameter for trading-off

KKT condition is given by:

$$a_n \geq 0, t_n y(x_n) - 1 + \xi_n \geq 0, a_n (t_n y(x_n) - 1 + \xi_n) = 0$$

$$\mu_n \geq 0, \xi_n \geq 0, \mu_n \xi_n = 0$$

Lagrangian is written by:

$$L(w, b, a) = \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \xi_n - \sum_{n=1}^N a_n \{ t_n y(x_n) - 1 + \xi_n \} - \sum_{n=1}^N \mu_n \xi_n$$

where $\mu_n \geq 0, a_n \geq 0$

then $\frac{\partial L}{\partial w} = 0, \frac{\partial L}{\partial b} = 0$, we obtain $\{w = \sum_{n=1}^N a_n t_n \phi(x_n)$

$$\frac{\partial L}{\partial \xi_n} = 0$$

$$\sum_{n=1}^N a_n t_n = 0, a_n = C - \mu_n$$

and dual Lagrange is: $\min \{ \hat{L}(a) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(x_n, x_m) \}$

subject to $\begin{cases} 0 \leq a_n \leq C \\ \sum_{n=1}^N a_n t_n = 0 \end{cases} \rightarrow$ control the quantity of support vectors.

Solution Interpretation:

① $a_n = 0 \Rightarrow$ nonsupport vector

② $0 < a_n < C$, then $\mu_n > 0$, then $\xi_n = 0$, this point lies on the margin.

③ if $a_n = C$, then $\mu_n = 0, \xi_n \leq 1$, or $\xi_n > 1$
 (classified) (misclassified)

To determine b , support vectors a_n satisfy $0 < a_n < C, \xi_n = 0, t_n y(x_n) = 1$.

then we have $\ln \left(\sum_{n \in S} a_n t_n k(x_n, x_n) + b \right) = 1$

$$\Rightarrow b = \frac{1}{N_M} \sum_{n \in M} \left(t_n - \sum_{m \in S} a_m t_m k(x_n, x_m) \right)$$

$\rightarrow M$ is a set with data points having $0 < a_n < C$.

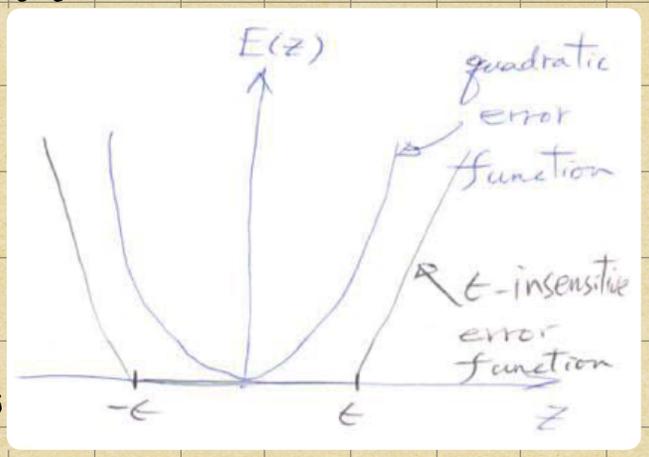
SVM for regression: (Support Vector Regression)

We define simple error function:

$$\Rightarrow \frac{1}{2} \sum_{n=1}^N (y_n - t_n)^2 + \frac{1}{2} \|w\|^2$$

To obtain sparse solution as:

$$E_\epsilon(y(x) - t) = \begin{cases} 0 & \text{if } |y(x) - t| \leq \epsilon \\ |y(x) - t| - \epsilon & \end{cases}$$



A new regularized error function:

$$C \sum_{n=1}^N E_\epsilon(y(x_n) - t_n) + \frac{1}{2} \|w\|^2$$

By introduce two slack variables:

$$\xi_n \geq 0 \equiv t_n > y(x_n) + \epsilon$$

$$\zeta_n \geq 0 \equiv t_n < y(x_n) - \epsilon$$

$$\text{For } y_n - \epsilon \leq t_n \leq y_n + \epsilon \Rightarrow \xi_n = \zeta_n = 0$$

Error function of SVR:

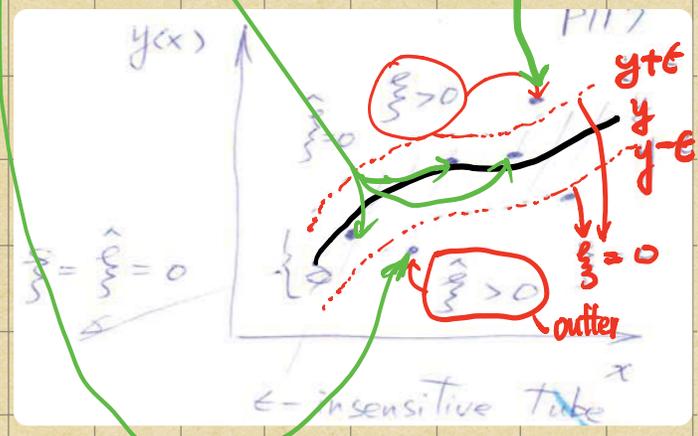
$$C \sum_{n=1}^N (\xi_n + \zeta_n) + \frac{1}{2} \|w\|^2$$

Constraints:

$$\xi_n \geq 0 \ \& \ \zeta_n \geq 0 \ \&$$

$$t_n \leq y(x_n) + \epsilon + \xi_n \ \&$$

$$t_n \geq y(x_n) - \epsilon - \zeta_n$$



Lagrange optimization:

$$L = C \sum_{n=1}^N (\xi_n + \zeta_n) + \frac{1}{2} \|w\|^2 - \sum_{n=1}^N (\alpha_n \xi_n + \hat{\alpha}_n \zeta_n) - \sum_{n=1}^N a_n (\epsilon + \xi_n + y_n - t_n) - \sum_{n=1}^N \hat{a}_n (t_n + \zeta_n - y_n + \epsilon)$$

$$\Rightarrow \begin{cases} \partial L / \partial w = 0 \Rightarrow w = \sum_{n=1}^N (\alpha_n - \hat{\alpha}_n) \phi(x_n) \\ \partial L / \partial b = 0 \Rightarrow \sum_{n=1}^N (\alpha_n - \hat{\alpha}_n) = 0 \end{cases}$$

$$\partial L / \partial \hat{\alpha}_n = 0 \Rightarrow \hat{\alpha}_n + \alpha_n = 0$$

$$\partial L / \partial \xi_n = 0 \Rightarrow \alpha_n \text{ and } \alpha_n = 0$$

Dual presentation:

$$\tilde{L}(a, \hat{a}) = -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N (a_n - \hat{a}_n)(a_m - \hat{a}_m) k(x_n, x_m) - t \sum_{n=1}^N (a_n + \hat{a}_n) + \sum_{n=1}^N (a_n - \hat{a}_n) t_n$$

From the result.

$$y(x) = \sum_{n=1}^N (a_n - \hat{a}_n) k(x, x_n) + b$$

The KKT conditions are given by:

$$a_n(\epsilon + \xi_n + y_n - t_n) = 0, \quad \hat{a}_n(\epsilon + \hat{\xi}_n - y_n + t_n) = 0$$

$$(c - a_n)\xi_n = 0, \quad (c - \hat{a}_n)\hat{\xi}_n = 0$$

The parameter "b" can be found by:

$$b = t_n - \epsilon - w^T \phi(x_n)$$

$$= t_n - t - \sum_{m=1}^N (a_m - \hat{a}_m) k(x_n, x_m)$$

7.2 Relance Vector Machine (RVM) [7.1 is all about SVM above & RVM is Sparse Kernel Machine]

Core idea:

- ① SVM is without probability, how about add Probability to it?
- ② SVM is two-class classification, how about more?
- ③ C or V should be found from held-out data.

The SVM-like form in RVM is: ^③ the number of w is much

② $y(x) = \sum_{n=1}^N w_n k(x, x_n) + b$, then likelihood function is larger than others

$P(t|X, w, \beta) = \prod_{n=1}^N P(t_n|x_n, w, \beta^{-1})$. And the prior distrib is $P(w|\alpha) = \prod_{i=1}^N N(w_i|0, \alpha_i^{-1})$ hyperparameter.

we could obtain posterior distrib: $P(w|t, X, \alpha, \beta) = N(w|m, \Sigma)$

where $m = \beta \Sigma \Phi^T t$ (w_{MAP}) $\Sigma = (C + \beta \Phi^T \Phi)^{-1}$

1. Estimate α, β (training)

α, β are determined by "evidence approximation":

$$P(t|X, \alpha, \beta) = \int P(t|X, w, \beta) P(w|\alpha) dw, \text{ then}$$

$$\ln P(t|X, \alpha, \beta) = \ln N(t|0, C) = -\frac{1}{2} \{ N \ln(2\pi) + \ln|C| + t^T C^{-1} t \}$$

where $C = \beta^{-1} I + \Phi A^T \Phi^T$. Let $\nabla \ln P(t|X, \alpha, \beta) = 0$

After amount of calculations we can obtain:

$$\alpha_i^{new} = r_i / m_i^2, \quad (\beta^{new})^{-1} = \frac{1 + \sum_i m_i^2}{N - \sum_i r_i}$$

$$P(t|x, X, \alpha^*, \beta^*) = \int P(t|x, w, \beta^*) \cdot P(w|X, t, \alpha^*, \beta^*) dw = N(t|m^T \phi(x), \sigma^2(x))$$

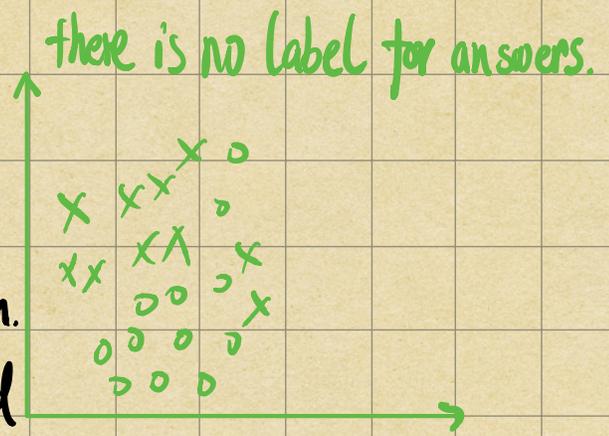
where $\sigma^2(x) = (\beta^*)^{-1} + \phi(x)^T \Sigma \phi(x)$

2. Predictive Distrib (test)

Day 7: Unsupervised Learning

clustering problem

- with probability: mixture of Gaussian.
- without probability: geometry-method



Chapter 9: Mixture Models and EM

E: expectation M: Maximization.

9.1 K-means Clustering (non-probability & multi-variables)

太简单, 不考.

Define $J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$ $r_{nk} = \begin{cases} 1 & \text{if } k = \operatorname{argmin}_j \|x_n - \mu_j\|^2 \\ 0 & \text{otherwise} \end{cases}$

↑ membership indicator

hard
soft

(the number of data is N & class is k)

31. then $\frac{\partial}{\partial \mu_k} \ln P(x) = 2 \sum_{n=1}^N \gamma_{nk} (x_n - \mu_k) = 0$

要考就考它!

9.2 Mixture of Gaussians: (with probability)

Given $p(x) = \sum_{k=1}^K \pi_k N(x | \mu_k, \Sigma_k)$
 \uparrow mixture weight.

Z is a latent variable represent one of the K latent states

$\therefore P(Z_k=1) = \pi_k$, then $P(x) = \sum_Z P(Z) \cdot P(x|Z) = \sum_{k=1}^K \pi_k N(x | \mu_k, \Sigma_k)$

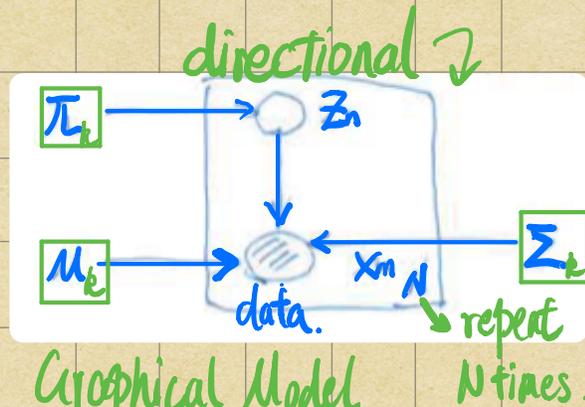
9.2.1 Maximum likelihood

Let $X = \{x_1, \dots, x_n\}$ & $Z = \{z_1, \dots, z_n\}$

likelihood function is:

$\ln P(x | \pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k N(x_n | \mu_k, \Sigma_k) \right\}$

$\sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k N(x_n | \mu_k, \Sigma_k) \right\}$



9.2.2: EM for Gaussian mixtures

$\frac{\partial}{\partial \mu_k} \ln P(x | \pi, \mu, \Sigma) = - \sum_{n=1}^N \frac{\pi_k N(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x_n | \mu_j, \Sigma_j)} \cdot \Sigma_k^{-1} (x_n - \mu_k) = 0$

$\Rightarrow \mu_k = \frac{1}{N_k} \sum_{n=1}^N r(z_{nk}) x_n$, where $N_k = \sum_{n=1}^N r(z_{nk})$

Let $\frac{\partial}{\partial \Sigma_k} \ln P(x | \pi, \mu, \Sigma) = 0 \Rightarrow \Sigma_k = \frac{1}{N_k} \sum_{n=1}^N r(z_{nk}) (x_n - \mu_k) \cdot (x_n - \mu_k)^T$

$\frac{\partial}{\partial \pi_k} (\ln P(x | \pi, \mu, \Sigma) + \lambda (\sum_{k=1}^K \pi_k - 1)) = \frac{N}{N_k} \frac{N(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x_n | \mu_j, \Sigma_j)} + \lambda$

Lagrange multiplier

$\lambda = -N \rightarrow \pi_k = N_k / N$

9.2.3: EM for Gaussian mixtures: (Algorithm)

① Initialize μ_k, Σ_k , evaluate log likelihood.

② E-step: $r(z_{nk}) = \frac{\pi_k N(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x_n | \mu_j, \Sigma_j)}$

③ M-step: $\mu_k^{new} = \frac{1}{N_k} \sum_{n=1}^N r(z_{nk}) \cdot x_n$

$\Sigma_k^{new} = \frac{1}{N_k} \sum_{n=1}^N r(z_{nk}) \cdot (x_n - \mu_k^{new}) \cdot (x_n - \mu_k^{new})^T$

32.

9.3 An alternative View of EM:

$$\ln P(X|\theta) = \ln \left\{ \sum_{\mathbf{z}} P(X, \mathbf{z}|\theta) \right\}$$

↑ incomplete and $\{X, \mathbf{z}\}$ is complete data.

General EM Algorithm:

$$E\text{-step: } Q(\theta, \theta^{\text{old}}) = E_{\mathbf{z}} [\ln P(X, \mathbf{z}|\theta) | X, \theta^{\text{old}}]$$

$$M\text{-step: } Q^{\text{new}} = \underset{\theta}{\text{argmax}} Q(\theta, \theta^{\text{old}})$$

Gaussian mixture revisited:

$$\ln P(X, \mathbf{z} | \mu, \Sigma, \tau) = \sum_{n=1}^N \sum_{k=1}^K \mathbb{1}_{\{z_{nk}\}} \{ \ln \tau_k + \ln N(X_n | \mu_k, \Sigma_k) \}$$

$$\therefore P(X, \mathbf{z}) = P(\mathbf{z}) P(X | \mathbf{z}) = \prod_k \tau_k^{\sum_n \mathbb{1}_{\{z_{nk}\}}} \prod_k N(X | \mu_k, \Sigma_k)^{\sum_n \mathbb{1}_{\{z_{nk}\}}}$$

$$E_{\mathbf{z}} [\ln P(X, \mathbf{z} | \mu, \Sigma, \tau)] = \sum_{n=1}^N \sum_{k=1}^K r(z_{nk}) \{ \ln \tau_k + \ln N(X_n | \mu_k, \Sigma_k) \}$$

$$\therefore E[z_{nk}] = \frac{\sum_{z_{nk}} z_{nk} [\tau_k N(X_n | \mu_k, \Sigma_k)]^{\sum_n z_{nk}}}{\sum_{z_{nj}} [\tau_j N(X_n | \mu_j, \Sigma_j)]^{\sum_n z_{nj}}} = r(z_{nk})$$

9.4: The EM Algorithm in general

$$P(X|\theta) = \sum_{\mathbf{z}} P(X, \mathbf{z}|\theta)$$

$$\begin{aligned} \ln P(X|\theta) &= \ln \sum_{\mathbf{z}} \frac{P(X, \mathbf{z}|\theta)}{q(\mathbf{z})} \cdot q(\mathbf{z}) = \ln E_{q(\mathbf{z})} \left[\frac{P(X, \mathbf{z}|\theta)}{q(\mathbf{z})} \right] \geq E_{q(\mathbf{z})} \left[\ln \frac{P(X, \mathbf{z}|\theta)}{q(\mathbf{z})} \right] \\ &= \mathcal{L}(q, \theta) = \sum_{\mathbf{z}} q(\mathbf{z}) \ln \left\{ \frac{P(X, \mathbf{z}|\theta)}{q(\mathbf{z})} \right\}. \end{aligned}$$

$$\ln P(X|\theta) = \mathcal{L}(q, \theta) + \text{KL}(q||p)$$

$$\text{where } \text{KL}(q||p) = - \sum_{\mathbf{z}} q(\mathbf{z}) \ln \left\{ \frac{P(\mathbf{z} | X, \theta)}{q(\mathbf{z})} \right\}$$

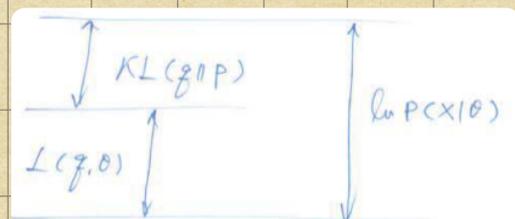


Illustration for EM algorithm:

$$\textcircled{1} \theta^{\text{old}}: \hat{q} = \underset{q}{\text{argmin}} \text{KL}(q||p) = 0, \hat{q}(\mathbf{z}) = P(\mathbf{z} | X, \theta^{\text{old}})$$

$$\textcircled{2} \hat{\theta}^{\text{new}}: \underset{\theta}{\text{argmax}} \mathcal{L}(\hat{q}, \theta)$$

$$\begin{aligned} \mathcal{L}(q, \theta) &= \sum_{\mathbf{z}} q(\mathbf{z}) \ln \left\{ \frac{P(X, \mathbf{z}|\theta)}{q(\mathbf{z})} \right\} = \sum_{\mathbf{z}} P(\mathbf{z} | X, \theta^{\text{old}}) \ln P(X, \mathbf{z}|\theta) - \\ &\quad \sum_{\mathbf{z}} P(\mathbf{z} | X, \theta^{\text{old}}) \ln P(\mathbf{z} | X, \theta^{\text{old}}) \\ &= \mathcal{L}(\theta, \theta^{\text{old}}) + \text{const} \end{aligned}$$

14 Information Theory: we define $h(x,y) = h(x) + h(y)$ & $h(x) = -\log_2 P(x)$

$$H[X] = -\sum_i P(x_i) \log_2 P(x_i), \quad H[X|Y] = -\sum_i P(x_i|y) \ln P(x_i|y) = E[-\ln P(X|Y)]$$

Maximum Entropy (ME): $\tilde{P} = -\sum P(x_i) \ln P(x_i) + \lambda (\sum P(x_i) - 1)$: Given X is continuous

random variables $\Rightarrow H[X] = -\int P(x) \ln P(x) dx$ (differential entropy)