

PRML 1. INTRODUCTION

- Supervised learning { classification ... discrete output
regression ... continuous output
- Unsupervised learning { clustering
density estimation
visualization etc.
- Reinforcement learning : exploration \leftrightarrow exploitation

§ 1.1

"regularization"

("shrinkage" or "weight decay")

$$\tilde{E}(w) = \frac{1}{2} \sum_n [y(x_n, w) - t_n]^2 + \frac{\lambda}{2} \|w\|^2$$

expectation
↓
data
↓
 $\sim 10^{-6}$

quadratic \Rightarrow "ridge regression"

const. term can be omitted.

§ 1.2

$$\text{Bayes' theorem } p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)} = \frac{p(X|Y)p(Y)}{\sum_Y p(X|Y)p(Y)}$$

$$\left(= \frac{p(X,Y)}{p(X)} \right)$$

for given params

Maximum likelihood : a frequentists' estimator $p(D|w)$
(or $-\ln p(D|w)$) ↑
the prob. to obtain data D

[cf. "bootstrap" method for uncertainty of w]

Gaussian distribution

$$N(\mathbf{X} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \exp \left[-\frac{1}{2} (\mathbf{X} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{X} - \boldsymbol{\mu}) \right]$$

$$\text{maximal likelihood } \boldsymbol{\mu}_{ML} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$$

$$\sigma_{ML}^2 = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}_{ML})^2$$

$$\text{but, } \mathbb{E}(\boldsymbol{\mu}_{ML}) = \boldsymbol{\mu}, \quad \mathbb{E}(\sigma_{ML}^2) = \frac{N-1}{N} \sigma^2 : \text{biased}$$

- Curve-fitting: predict $t = y(x, \mathbf{w}) = \sum_{i=0}^M w_i x^i$

from noised-data $\{\mathbf{x}, t\}$ assuming Gaussian noise β

$$\left[\text{i.e. } p(t | \mathbf{x}, \mathbf{w}, \beta) = N(t | y(\mathbf{x}, \mathbf{w}), \beta') \right]$$

\Rightarrow likelihood function

$$\ln p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum [y(\mathbf{x}_i, \mathbf{w}) - t_i]^2 + \frac{N}{2} \ln \beta + \text{const.}$$

$$\rightarrow \mathbf{w}_{ML} \text{ and } \beta_{ML} \text{ are obtained } \frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{n=1}^N \{y(\mathbf{x}_n, \mathbf{w}_{ML}) - t_n\}^2.$$

We can introduce a prior. e.g. $p(\mathbf{w} | \alpha) = N(\mathbf{w} | \mathbf{0}, \alpha^{-1} \mathbf{I}) = \left(\frac{\alpha}{2\pi} \right)^{(M+1)/2} \exp \left\{ -\frac{\alpha}{2} \mathbf{w}^T \mathbf{w} \right\}$

\Rightarrow likelihood hyperparameter: parameter controlling the models.

$$p(\mathbf{w} | \mathbf{x}, \mathbf{t}, \beta, \alpha) = \frac{p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) p(\mathbf{w} | \alpha)}{\sum_{\mathbf{w}} \text{(numerator)}} \quad \dots (2)$$

$$\ln p(\mathbf{w} | \mathbf{x}, \mathbf{t}, \beta, \alpha) = (1) - \frac{\alpha}{2} \|\mathbf{w}\|^2 + \text{const.} \quad \text{new information}$$

$\dots \mathbf{w}_{ML}$ is given by "regularized" $\tilde{\mathbf{E}}$ in §1.1.

with $\lambda = \alpha/\beta$.

(but this is "point estimation" \dots not Bayesian way)

- Curve-fitting in Bayesian way

Bayesian approach: consider $p(t|x, \mathbf{x}, \mathbf{t})$

$$p(t|x, \mathbf{x}, \mathbf{t}) = \int p(t|x, w) p(w|\mathbf{x}, \mathbf{t}) dw$$

we knew ② ③ [α and β are assumed fixed]

$$p(t|x, \mathbf{x}, \mathbf{t}) = \mathcal{N}(t|m(x), s^2(x)) \quad (1.69)$$

where the mean and variance are given by



$$m(x) = \beta \phi(x)^T \mathbf{S} \sum_{n=1}^N \phi(x_n) t_n \quad (1.70)$$

$$s^2(x) = \beta^{-1} + \phi(x)^T \mathbf{S} \phi(x). \quad (1.71)$$

Here the matrix \mathbf{S} is given by we had in previous approach new contribution to the uncertainty.

$$\mathbf{S}^{-1} = \alpha \mathbf{I} + \beta \sum_{n=1}^N \phi(x_n) \phi(x_n)^T \quad (1.72)$$

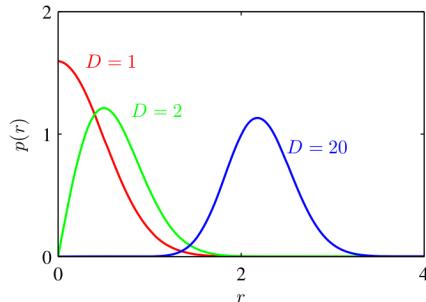
where \mathbf{I} is the unit matrix, and we have defined the vector $\phi(x)$ with elements $\phi_i(x) = x^i$ for $i = 0, \dots, M$.

§ 1.3 Model Selection

- separate data into training + validation (+ test) set
- cross validation: S -times run w. $\frac{S-1}{S}$ training + 1 validation
 (leave-one-out for $S=N$)
 - heavy calculation
 - difficult for multi-parameter model selection
- Akaike information criterion $\ln p(D| \text{data}) - M$) prefers overly-simple models
- Bayesian information criterion

§ 1.4 Curse of dimensionality

Figure 1.23 Plot of the probability density with respect to radius r of a Gaussian distribution for various values of the dimensionality D . In a high-dimensional space, most of the probability mass of a Gaussian is located within a thin shell at a specific radius.



§ 1.5 Decision theory

$$p(C_k | \mathcal{X}) = \frac{p(\mathcal{X} | C_k) p(C_k)}{p(\mathcal{X})}$$

↑ prior

"class" ↑ observed data ↑

- to assign \mathcal{X} to a class
- \Leftrightarrow to separate the space of \mathcal{X} into regions R_k

$$p(\text{correct}) = \sum_k \int_{R_k} d\mathcal{X} p(\mathcal{X}, C_k)$$

• to maximise $p(\text{correct})$: k for \mathcal{X} is what gives the largest $p(\mathcal{X}, C_k)$

$$\Leftrightarrow \quad \quad \quad p(C_k | \mathcal{X})$$

• to minimize "expected loss" $L_{k,j}$

↑ true class ↓ assigned class

$$\mathbb{E}(L) = \sum_{k,j} \int_{R_j} L_{kj} p(\mathcal{X}, C_k) d\mathcal{X}$$

\Leftrightarrow choose j that minimizes $\sum_k L_{kj} p(\mathcal{X}, C_k)$

$$\left(\Leftrightarrow \sum_k L_{kj} p(C_k | \mathcal{X}) \right)$$

• "reject" option (postpone our decision)

$$\max_k p(C_k | \mathcal{X}) < \theta$$

(Note that this is a good criterion.
Also θ should be $\in (\frac{1}{k}, 1)$)

① Classification problem = Inference + Decision.

(a) calculate $p(C_k, \mathcal{X}), p(\mathcal{X})$, etc. and then decide

• too much computation ↑ very useful to detect rare cases [outlier detection, novelty detection]

(b) calculate $p(C_k | \mathcal{X})$ and then decide.

(c) find a function $f: \mathcal{X} \mapsto k$

• no access to posterior probabilities.

Decision theory for Regression

$$\mathbb{E}[L] = \int d\mathbf{x} dt L(t, \underbrace{y(\mathbf{x})}_{\substack{\text{estimate} \\ \mathbb{R}^N \rightarrow \mathbb{R}}}) p(\mathbf{x}, t)$$

↑
input
 $\in \mathbb{R}^N$

↑
true
 $\in \mathbb{R}$

$$\text{for } L = [t - y(\mathbf{x})]^2,$$

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

$\therefore \hat{y}(\mathbf{x}) = \frac{\int dt t p(\mathbf{x}, t)}{p(\mathbf{x})} = \int dt t p(t|\mathbf{x}) \equiv \mathbb{E}_t[t|\mathbf{x}]$

dep. on \mathbf{x}
indep. of t for averaged t under the condition \mathbf{x}

$$\begin{aligned} [\text{OR: } \mathbb{E}[L]] &= \int d\mathbf{x} dt \left\{ [y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]]^2 p(\mathbf{x}, t) \right. \\ &\quad + 2(y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]) (\mathbb{E}[t|\mathbf{x}] - t) p(\mathbf{x}, t) \\ &\quad \left. + [\mathbb{E}[t|\mathbf{x}] - t]^2 p(\mathbf{x}, t) \right\} \xrightarrow{\text{cancel terms}} \mathbb{E}[p(\mathbf{x})] \mathbb{E}[t|\mathbf{x}] p(\mathbf{x}) \\ &= \int d\mathbf{x} [y(\mathbf{x}) - \mathbb{E}]^2 p(\mathbf{x}) + \int d\mathbf{x} dt [\mathbb{E}[t|\mathbf{x}] - t]^2 p(\mathbf{x}, t) \\ &\quad \hookrightarrow \int d\mathbf{x} \text{Var}_t[t|\mathbf{x}] p(\mathbf{x}) \end{aligned}$$

$$\cdot \mathbb{E}_t[f(t)|\mathbf{x}] = \int dt f(t) p(t|\mathbf{x})$$

$$\cdot \mathbb{V}[Y|\mathbf{x}] = \mathbb{E}\left[(Y - \mathbb{E}[Y|\mathbf{x}])^2 \mid \mathbf{x}\right]$$

$$\begin{aligned} \cdot \mathbb{V}_t[f(t)|\mathbf{x}] &= \mathbb{E}_t\left[(f(t) - \mathbb{E}[f(t)|\mathbf{x}])^2 \mid \mathbf{x}\right] \\ &= \int dt (f(t) - \mathbb{E}[f(t)|\mathbf{x}])^2 p(t|\mathbf{x}) \end{aligned}$$

§ 1.6 Entropy

discrete: $H = - \sum_{i=1}^N p_i \ln p_i \in [0, \ln N]$

continuous: $H[X] = - \int p(x) \ln p(x) dx$

$$\left[\begin{aligned} H_{\text{disc.}} &= - \sum p_i \ln p_i \\ &\rightarrow - \sum p(x) \Delta x \ln p(x) \Delta x \\ &= - \sum p(x) \ln p(x) \Delta x - \sum p(x) \ln \Delta x \Delta x \\ &= - \int p(x) \ln p(x) dx - \lim_{\Delta \rightarrow 0} \ln \Delta \\ &= H[X] + \infty \end{aligned} \right]$$

for mean μ and variance σ^2 ,

Gaussian distrib. gives maximal entropy $\frac{1 + \log 2\pi\sigma^2}{2}$

- Conditional entropy $H[Y|X] = - \int p(x, y) \ln P(Y|X) dx dy$
 $= H[X, Y] - H[X]$

- Relative entropy (Kullback-Leibler divergence)

$$KL(P||Q) = - \int p(x) \ln \frac{q(x)}{p(x)} dx$$

where we modeled the true distrib. $p(x)$ by an approx. one $q(x)$.

Since we optimize our code by $q(x)$,

we have to pay extra bandwidth to send the information of outputs from P .
 $KL(P||Q)$

- $KL(P||Q) > 0$ for $P \neq Q$

- $KL(P||Q) = 0$ for $P = Q$

- Mutual information

$$\begin{aligned} I[X, Y] &= KL(P(X, Y) || P(X)P(Y)) \\ &= - \int p(x, y) \ln \frac{p(x)p(y)}{p(x, y)} dx dy \quad (\geq 0) \end{aligned}$$

$$H[X] = H[X|Y] + I[X, Y]$$

PRML 2. PROBABILITY DISTRIBUTIONS

§ 2.1 Binary Variables

- $\text{Bern}(x|\mu) = \mu^x (1-\mu)^{1-x}$ $[x=0,1]$ $E = \mu, V = \mu(1-\mu)$

- $\text{Bin}(x|N,\mu) = {}_N C_x \mu^x (1-\mu)^{N-x}$ $E = N\mu, V = N\mu(1-\mu)$

- $\text{Beta}(x|\alpha, \beta) = \frac{P(\alpha+b)}{P(\alpha)P(b)} x^{\alpha-1} (1-x)^{\beta-1}$ $E = \frac{\alpha}{\alpha+\beta}, V = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$
 $\stackrel{\uparrow}{[0,1]}$ $\underbrace{\hspace{10em}}$ $\underbrace{\hspace{10em}}_{\sim \alpha^{-1} \text{ for } \alpha=\beta}$
 useful for prior

Using Beta($x|\alpha, \beta$) as a prior, posterior after "1xm + 0xl" obs. is

$$\begin{aligned} p(x|m, l, \alpha, \beta) &\propto \mu^{m+\alpha-1} (1-\mu)^{l+\beta-1} \\ &= \text{Beta}(x|\alpha+m, \beta+l) \end{aligned}$$

"sequential learning"

General discussion on Variance in sequential learning

$$\mathbb{E}_\theta[\theta] = \mathbb{E}_{\mathcal{D}} [\mathbb{E}_\theta[\theta|\mathcal{D}]]$$

\uparrow prior mean of a parameter $\overbrace{\hspace{10em}}$ averaged over data possibility
 $\underbrace{\hspace{10em}}$ posterior mean of θ after \mathcal{D} -observation

(2.21)

where

$$\mathbb{E}_\theta[\theta] \equiv \int p(\theta) \theta d\theta \quad (2.22)$$

$$\mathbb{E}_{\mathcal{D}}[\mathbb{E}_\theta[\theta|\mathcal{D}]] \equiv \int \left\{ \int \theta p(\theta|\mathcal{D}) d\theta \right\} p(\mathcal{D}) d\mathcal{D} \quad (2.23)$$

says that the posterior mean of θ , averaged over the distribution generating the data, is equal to the prior mean of θ . Similarly, we can show that

$$\text{var}_\theta[\theta] = \mathbb{E}_{\mathcal{D}} [\text{var}_\theta[\theta|\mathcal{D}]] + \text{var}_{\mathcal{D}} [\mathbb{E}_\theta[\theta|\mathcal{D}]]. \quad (2.24)$$

prior variance = "averaged" posterior variance + (positive)

§ 2.2 Multinomial variables

Observation of "i" among "K" states $\Rightarrow \mathbf{x} = \{0, \dots, 0, 1, 0, \dots, 0\}$

If the prob. to get "i" is μ_i ,

$$\sum \mu_i = 1, \quad p(\mathbf{x}|\boldsymbol{\mu}) = \prod \mu_i^{x_i} \quad \cdots \sum x_i = 1$$

$$\mathbb{E}[\mathbf{x}|\boldsymbol{\mu}] = \sum_{\mathbf{x}} p(\mathbf{x}|\boldsymbol{\mu}) \mathbf{x} = \{p(x_1|\boldsymbol{\mu}), \dots\} = \boldsymbol{\mu}^T$$

Likelihood to observe $D = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$

$$p(D|\boldsymbol{\mu}) = \prod_{n=1}^N p(\mathbf{x}_n|\boldsymbol{\mu}) = \prod_{n=1}^N \prod_{k=1}^K \mu_k^{(x_n)_k} = \prod_k \mu_k^{\sum_n (x_n)_k}$$

$$= \prod_k \mu_k^{m_k} \quad \begin{matrix} m_k \leftarrow \# \text{ of observation of data } k \\ \dots \text{ "sufficient statistics"} \end{matrix}$$

$$\Rightarrow (\mu_k)_{ML} = m_k / N$$

$$\text{• Mult}(m_1, \dots, m_K | \boldsymbol{\mu}, N) = \frac{N!}{m_1! \dots m_N!} \prod_{k=1}^K \mu_k^{m_k} \quad (\sum m_k = N)$$

$$\text{• Dir}(\boldsymbol{\mu} | \boldsymbol{\alpha}) = \frac{P(\alpha_0)}{P(\alpha_1) \dots P(\alpha_K)} \prod_{k=1}^K \mu_k^{\alpha_k - 1}$$

$$\text{where } \alpha_0 = \sum \alpha_k, \quad \sum \mu_k = 1.$$

prior $\text{Dir}(\boldsymbol{\mu} | \boldsymbol{\alpha}) = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_K)} \prod_{k=1}^K \mu_k^{\alpha_k - 1} \quad (2.38)$

posterior $p(\boldsymbol{\mu} | \mathcal{D}, \boldsymbol{\alpha}) = \text{Dir}(\boldsymbol{\mu} | \boldsymbol{\alpha} + \mathbf{m})$

$$= \frac{\Gamma(\alpha_0 + N)}{\Gamma(\alpha_1 + m_1) \dots \Gamma(\alpha_K + m_K)} \prod_{k=1}^K \mu_k^{\alpha_k + m_k - 1} \quad (2.41)$$

§ 2.3 Gaussian Distribution

$$N(\mathbf{x} | \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \exp\left[-\frac{1}{2}\Delta^2\right] \cdot \text{Unimodal}$$

... too flexible
... too restricted

$$\text{w. } \Delta^2 = (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)$$

- Σ : Real Hermitian (symmetric) & positive definite

\Rightarrow · diagonalizable (eigenvalues $\in \mathbb{R}$)

$$\Sigma = \mathbf{U}^T \mathbf{L} \mathbf{U}, \quad \Delta^2 = (\mathbf{x} - \mu)^T \mathbf{U} \mathbf{L}^{-1} \mathbf{U} (\mathbf{x} - \mu) =: \mathbf{y}^T \mathbf{L}^{-1} \mathbf{y}$$

$$\cdots N(\mathbf{x} | \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \exp\left[-\frac{\mathbf{y}^2}{2\lambda_i}\right]$$

- Writing N in terms of $\mathbf{y} = \mathbf{U}(\mathbf{x} - \mu)$

$$\begin{aligned} p(\mathbf{y}) &= \left| \frac{d\mathbf{x}}{d\mathbf{y}} \right| p(\mathbf{x}) = |\mathbf{U}|^{-1} p(\mathbf{x}) \\ &= \prod_{i=1}^D \frac{1}{\sqrt{2\pi\lambda_i}} \exp\left(-\frac{y_i^2}{2\lambda_i}\right) \end{aligned} \quad \begin{cases} |\mathbf{U}| = 1 \\ |\Sigma| = |\mathbf{L}| = \prod \lambda_i \end{cases}$$

- $E(\mathbf{x}) = \mu$

- $Cov[\mathbf{x}] = \Sigma$

for vector variables.

$$\begin{aligned} Cov[\mathbf{x}, \mathbf{y}] &= E((\mathbf{x} - E(\mathbf{x}))(\mathbf{y} - E(\mathbf{y}))^T) \\ &= E(\mathbf{x}\mathbf{y}^T) - E(\mathbf{x})E(\mathbf{y}^T) \end{aligned}$$

- Maximal Likelihood $p(\mathbf{x} | \mu, \Sigma)$

sufficient statistics

$$\sum \mathbf{x}_n, \sum \mathbf{x}_n \mathbf{x}_n^T \Rightarrow \mu_{MC} = \frac{1}{N} \sum \mathbf{x}_n, \quad \Sigma_{MC} = \frac{1}{N} \sum (\mathbf{x}_n - \mu_{MC})(\mathbf{x}_n - \mu_{MC})^T$$

$$\left(\text{again, } E(\mu_{MC}) = \mu \text{ but } E(\Sigma_{MC}) = \frac{N-1}{N} \Sigma \right)$$

- $P(\mathbf{x}_a, \mathbf{x}_b)$ is Gaussian \implies $\begin{cases} P(\mathbf{x}_a | \mathbf{x}_b) \\ P(\mathbf{x}_b) \end{cases}$ are Gaussian.

§2.3.1 Conditional Gaussian distribution

$$p(\mathbf{x} | \mathbf{y}) = \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{y})} \xrightarrow{\text{fixed}} \text{normalization}$$

Separating \mathbf{x} into $(\mathbf{x}_a, \mathbf{x}_b)$ and $\Lambda = \Sigma^{-1}$ into $\begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ab}^T & \Lambda_{bb} \end{pmatrix}$,

$$\Delta^2 = (\mathbf{x}_a - \mu_a)^T \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ab}^T & \Lambda_{bb} \end{pmatrix} (\mathbf{x}_a - \mu_a)$$

$$= (\mathbf{x}_a - \mu_{a|b})^T \Lambda_{ab} (\mathbf{x}_a - \mu_{a|b}) + (\text{Remainders}) \xrightarrow{\text{renormalized}}$$

$$\begin{aligned} \text{w.t. } \mu_{a|b} &= \mu_a - \Lambda_{aa}^{-1} \Lambda_{ab} (\mathbf{x}_b - \mu_b) \\ &= \mu_a + \Sigma_{ab} \Sigma_{bb}^{-1} (\mathbf{x}_b - \mu_b). \end{aligned} \quad \left. \begin{aligned} (\mathbf{A}\mathbf{B})^{-1} &= \begin{pmatrix} \mathbf{M} & -\mathbf{MBD}^{-1} \\ -\mathbf{D}^T \mathbf{C} \mathbf{M} & \mathbf{D}^{-1} + \mathbf{D}^T \mathbf{C} \mathbf{M} \mathbf{B} \mathbf{D}^{-1} \end{pmatrix} \\ \text{w.t. } \mathbf{M} &= (\mathbf{A} - \mathbf{B} \mathbf{D}^T \mathbf{C})^{-1} \end{aligned} \right\} \begin{aligned} \Lambda_{aa} &= (\Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ab}^T)^{-1} \\ \Lambda_{ab} &= -\Lambda_{aa} \Sigma_{ab} \Sigma_{bb}^{-1} \\ \Lambda_{bb} &= \Sigma_{bb}^{-1} + \Sigma_{bb}^{-1} \Sigma_{ab}^T \Lambda_{aa} \Sigma_{ab} \Sigma_{bb}^{-1} \end{aligned}$$

$$\begin{aligned} \Sigma_{ab} &= \Lambda_{aa}^{-1} \\ &= \Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ab}^T \end{aligned}$$

§2.3.2 Marginal Gaussian distribution

$$p(\mathbf{x}_a) = \int p(\mathbf{x}_a, \mathbf{x}_b) d\mathbf{x}_b$$

$$\boxed{\mathbb{E}[\mathbf{x}_a] = \mu_a}$$

$$\boxed{\text{Cov}[\mathbf{x}_a] = \Sigma_{aa}}$$

Partitioned Gaussians

Given a joint Gaussian distribution $\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \Sigma)$ with $\Lambda \equiv \Sigma^{-1}$ and

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix}, \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix} \quad (2.94)$$

$$\Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}, \quad \Lambda = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix}. \quad (2.95)$$

Conditional distribution:

$$p(\mathbf{x}_a | \mathbf{x}_b) = \mathcal{N}(\mathbf{x}_a | \boldsymbol{\mu}_{a|b}, \Lambda_{aa}^{-1}) \quad (2.96)$$

$$\boldsymbol{\mu}_{a|b} = \boldsymbol{\mu}_a - \Lambda_{aa}^{-1} \Lambda_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b). \quad (2.97)$$

Marginal distribution:

$$p(\mathbf{x}_a) = \mathcal{N}(\mathbf{x}_a | \boldsymbol{\mu}_a, \Sigma_{aa}). \quad (2.98)$$

§ 2.3.3 Bayes' theorem

$$\begin{cases} p(x) = \mathcal{N}(x | \mu_x, \Lambda^{-1}) \\ p(y|x) = \mathcal{N}(y | Ax + b, L^{-1}) \end{cases} \xrightarrow{\text{construct}} \begin{cases} p(x|y) \\ p(y) \end{cases}$$

linear in x indep. of x "linear Gaussian model"

Marginal and Conditional Gaussians

Given a marginal Gaussian distribution for x and a conditional Gaussian distribution for y given x in the form

$$p(x) = \mathcal{N}(x | \mu, \Lambda^{-1}) \quad (2.113)$$

$$p(y|x) = \mathcal{N}(y | Ax + b, L^{-1}) \quad (2.114)$$

the marginal distribution of y and the conditional distribution of x given y are given by

$$p(y) = \mathcal{N}(y | A\mu + b, L^{-1} + \Lambda\Lambda^{-1}A^T) \quad (2.115)$$

$$p(x|y) = \mathcal{N}(x | \Sigma \{ A^T L(y - b) + \Lambda\mu \}, \Sigma) \quad (2.116)$$

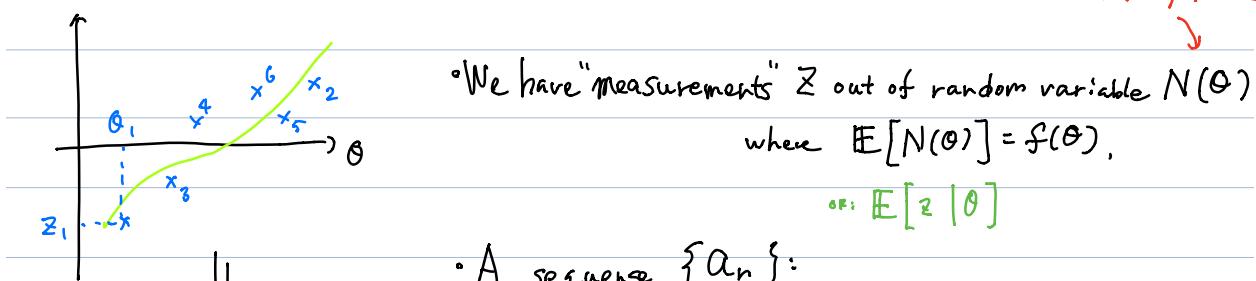
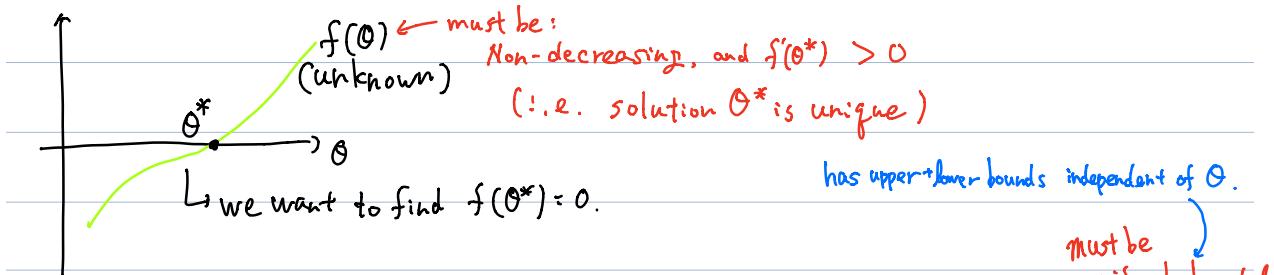
where

$$\Sigma = (\Lambda + A^T L A)^{-1}. \quad (2.117)$$

(Straight forward w. $p(x) \equiv p(x,y) = p(x)p(y|x)$)

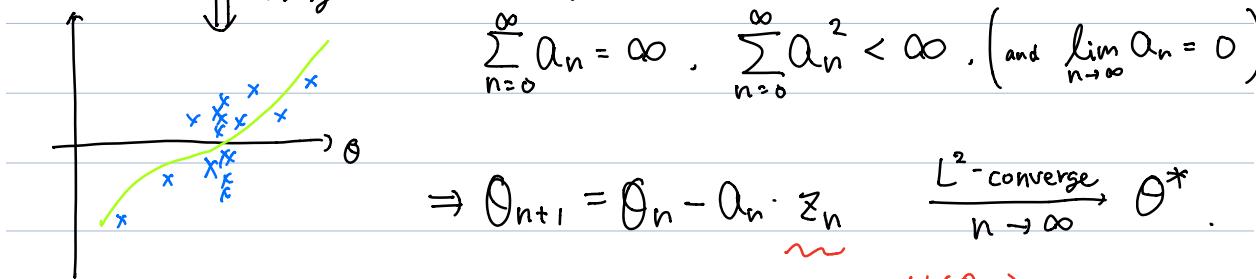
* § 2.3.5

★ Robbins-Monro algorithm



A sequence $\{\alpha_n\}$:

$$\sum_{n=0}^{\infty} \alpha_n = \infty, \quad \sum_{n=0}^{\infty} \alpha_n^2 < \infty, \quad (\text{and } \lim_{n \rightarrow \infty} \alpha_n = 0)$$



II Sequential ML Estimation

μ_{ML} maximizes $\prod_{n=1}^N p(x_n | \mu)$: e. solves $-\frac{1}{N} \frac{\partial}{\partial \mu} \sum_{n=1}^N \ln p(x_n | \mu) = 0$.

Thus, the ideal value of μ_{ML} solves $\mathbb{E} \left[\underbrace{-\frac{\partial}{\partial \mu} \ln p(x | \mu)}_{N(\theta)} \right] \Big|_{\mu=\mu_{ML}} = 0$

$$\therefore \mu_{ML}^{(N+1)} = \mu_{ML}^{(N)} - \alpha_n \cdot \left(-\frac{\partial}{\partial \mu} \ln p(x_n | \mu) \Big|_{\mu=\mu_{ML}^{(N)}} \right)$$

will L^2 -converge (almost sure converge) to $(\mu_{ML})^{\text{ideal}}$

$$\because \text{for Gaussian, } \mu_{ML}^{(N+1)} = \mu_{ML}^{(N)} + \alpha_n \cdot \frac{x_n - \mu_{ML}^{(N)}}{\sigma^2}$$

§2.3.6 Bayesian treatment

Remember: $p(\theta|x) \propto p(x|\theta)p(\theta)$

- Unknown μ , known σ^2

$$p(\mathbf{x}|\mu) \propto \exp\left[-\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2\right] \quad \text{posterior = prior} \oplus \text{data}$$

$$\Rightarrow \text{A good prior: } N(\mu| \mu_0, \sigma_0^2) \quad \begin{aligned} \mu_N &= \frac{\sigma^2}{N\sigma_0^2 + \sigma^2}\mu_0 + \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2}\mu_{ML} \\ \frac{1}{\sigma_N^2} &= \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2} \end{aligned} \quad (2.141) \quad (2.142)$$

- Known μ , unknown $\sigma^2 \equiv \lambda^{-1/2}$

$$p(\mathbf{x}|\lambda) \propto \lambda^{N/2} \exp\left[-\frac{\lambda}{2} \sum_{n=1}^N (x_n - \mu)^2\right]$$

$$\text{Gam}(\lambda|a,b) = \frac{b^a}{\Gamma(a)} \lambda^{a-1} e^{-b\lambda} \quad \begin{aligned} a_N &= a_0 + \frac{N}{2} \\ b_N &= b_0 + \frac{1}{2} \sum_{n=1}^N (x_n - \mu)^2 = b_0 + \frac{N}{2} \sigma_{ML}^2 \end{aligned} \quad (2.150) \quad (2.151)$$

- Both Unknown

$$p(\mathbf{x}|\mu, \lambda) \propto \lambda^{N/2} \exp\left[-\frac{\lambda}{2} \sum (x_n - \mu)^2\right]$$

$$= \lambda^{N/2} \exp\left[-\frac{N\lambda}{2} ((\mu - \mu_{ML})^2 + \sigma_{ML}^2)\right]$$

$$\propto \left[\lambda^{1/2} \exp\left(-\frac{\lambda\mu^2}{2}\right)\right]^N \exp\left\{\lambda\mu \sum_{n=1}^N x_n - \frac{\lambda}{2} \sum_{n=1}^N x_n^2\right\}. \quad (2.152)$$

Gaussian-gamma distribution

$$GG(\mu, \lambda | \mu_0, \beta, a, b) = N(\mu | \mu_0, (\beta\lambda)^{-1}) \text{Gam}(\lambda | a, b)$$

* Multivariable case

μ -unknown: Gaussian

$$\Sigma\text{-unknown: Wishart} \quad W(\Lambda | \mathbf{W}, \nu) = B|\Lambda|^{(\nu-D-1)/2} \exp\left(-\frac{1}{2}\text{Tr}(\mathbf{W}^{-1}\Lambda)\right) \quad (2.155)$$

both unknown:

$$p(\mu, \Lambda | \mu_0, \beta, \mathbf{W}, \nu) = N(\mu | \mu_0, (\beta\Lambda)^{-1}) W(\Lambda | \mathbf{W}, \nu) \quad (2.157)$$

where $B(\mathbf{W}, \nu) = |\mathbf{W}|^{-\nu/2} \left(2^{\nu D/2} \pi^{D(D-1)/4} \prod_{i=1}^D \Gamma\left(\frac{\nu+1-i}{2}\right)\right)^{-1}. \quad (2.156)$

§2.3.7 Student's t-distribution : infinite mixture of Gaussian

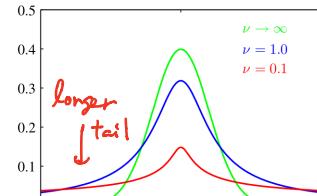
$$St(x|\mu, \lambda, \nu) = \frac{\Gamma(\frac{\nu+1}{2})}{\Gamma(\nu/2)} \left(\frac{\lambda}{\pi\nu} \right)^{\nu/2} \left[1 + \frac{\lambda(x-\mu)^2}{\nu} \right]^{-(\nu+1)/2}$$

this appears in the marginal distrib.

$$\begin{aligned} p(\mu|\mu_0, \beta=1, a, b) &= \int_0^\infty GG(\mu, \lambda|\mu_0, 1, a, b) d\lambda \\ &= \int_0^\infty N(\mu|\mu_0, \lambda^{-1}) \text{Gam}(\lambda|a, b) d\lambda \\ &= \frac{\Gamma(a+\frac{1}{2})}{\Gamma(a)} \left(\frac{1}{2\pi b} \right)^{1/2} \left[1 + \frac{(\mu-\mu_0)^2}{2b} \right]^{-\frac{2a+1}{2}} \end{aligned}$$

$$St(x|\mu, \lambda, \nu \rightarrow \infty) = N(x|\mu, \lambda^{-1})$$

$$St(x|\mu, \lambda, \nu = 1) = \text{Cauchy distribution}$$



↳ robustness

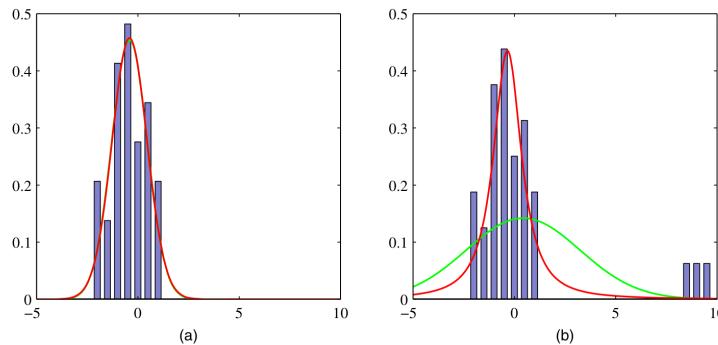


Figure 2.16 Illustration of the robustness of Student's t-distribution compared to a Gaussian. (a) Histogram distribution of 30 data points drawn from a Gaussian distribution, together with the maximum likelihood fit obtained from a t-distribution (red curve) and a Gaussian (green curve, largely hidden by the red curve). Because the t-distribution contains the Gaussian as a special case it gives almost the same solution as the Gaussian. (b) The same data set but with three additional outlying data points showing how the Gaussian (green curve) is strongly distorted by the outliers, whereas the t-distribution (red curve) is relatively unaffected.

* Multivariate version

$$St(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}, \nu) = \int_0^\infty N(\mathbf{x}|\boldsymbol{\mu}, (\eta\boldsymbol{\Lambda})^{-1}) \text{Gam}(\eta|\nu/2, \nu/2) d\eta. \quad (2.161)$$

$$= \frac{\Gamma(D/2 + \nu/2)}{\Gamma(\nu/2)} \frac{|\boldsymbol{\Lambda}|^{1/2}}{(\pi\nu)^{D/2}} \left[1 + \frac{\Delta^2}{\nu} \right]^{-D/2 - \nu/2} \quad (2.162) \quad \Delta^2 = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Lambda} (\mathbf{x} - \boldsymbol{\mu})$$

$$\mathbb{E}[\mathbf{x}] = \boldsymbol{\mu}, \quad \text{if } \nu > 1 \quad (2.164)$$

$$\text{cov}[\mathbf{x}] = \frac{\nu}{(\nu - 2)} \boldsymbol{\Lambda}^{-1}, \quad \text{if } \nu > 2 \quad (2.165)$$

$$\text{mode}[\mathbf{x}] = \boldsymbol{\mu} \quad (2.166)$$

§ 2.3.8 Periodic variable: von Mises distribution

$$P(\theta | \theta_0, m) = \frac{1}{2\pi I_0(m)} e^{m \cos(\theta - \theta_0)}$$

$\text{where } I_0(m) = \frac{1}{2\pi} \int_0^{2\pi} \exp\{m \cos \theta\} d\theta.$

$\tilde{[0, 2\pi]}$ periodic

this corresponds to $N\left(\begin{pmatrix} x \\ y \end{pmatrix} \middle| \begin{pmatrix} x_0 \\ y_0 \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}\right)$ restricted on $x^2 + y^2 = 1$.

$$\begin{cases} \theta_0^{\text{ML}} = \text{Arctan} \frac{\sum \sin \theta_n}{\sum \cos \theta_n} \\ \frac{I_1(m_{\text{ML}})}{I_0(m_{\text{ML}})} = \frac{1}{N} \sum \cos(\theta_n - \theta_0^{\text{ML}}) \\ = \left(\frac{1}{N} \sum \cos \theta_n \right) \cos \theta_0^{\text{ML}} + \left(\frac{1}{N} \sum \sin \theta_n \right) \sin \theta_0^{\text{ML}} \end{cases}$$

§ 2.3.9 Mixtures of Gaussians

$$p(x) = \sum_{k=1}^K \pi_k N(x | \mu_k, \Sigma_k) \quad \text{w. } \begin{cases} \pi_k \in [0, 1] \\ \sum \pi_k = 1. \end{cases}$$

We can view this as $P(x) = \sum_k p(k) P(x | k)$

Then

$$p(k|x) = \frac{p(k)p(x|k)}{\sum p(i)p(x|i)} = \frac{\pi_k N(x | \mu_k, \Sigma_k)}{\sum \pi_i N(x | \mu_i, \Sigma_i)}$$

$\equiv \gamma_k(x)$: responsibilities

§ 2.4 Exponential Family

$$p(x|\eta) = \underbrace{g(\eta)}_{\text{normalization}} \cdot h(x) e^{\eta^T \cdot u(x)}$$

$$\left(\equiv \exp \left[\eta^T \cdot u(x) + \tilde{g}(\eta) + \tilde{h}(x) \right] \right)$$

$$\text{Bern}(x|\mu) \propto e^{x \log \frac{\mu}{1-\mu}} : \begin{matrix} h(x) \\ 1 \\ \ln \frac{\mu}{1-\mu} \\ x \end{matrix}$$

$$\text{Beta}(x|a,b) \propto \frac{1}{x(1-x)} e^{a \log x + b \log(1-x)} : \begin{matrix} 1 \\ \frac{1}{x(1-x)} \\ a \\ b \\ \log x \\ \log(1-x) \end{matrix}$$

$$\text{Bin}(x|N,\mu) \propto \binom{N}{x} e^{x \log \frac{\mu}{1-\mu}} : \begin{matrix} N \\ C_x \\ \ln \frac{\mu}{1-\mu} \\ x \end{matrix}$$

$$\mathcal{N}(x|\mu, \lambda^{-1}) \propto e^{-\frac{\lambda}{2}(x^2 - 2x\mu)} : \begin{matrix} 1 \\ -\lambda/2 \\ \lambda \mu \\ x^2 \\ x \end{matrix}$$

$$\text{Pois}(n|\lambda) = \frac{\lambda^n e^{-\lambda}}{n!} : \begin{matrix} 1 \\ \frac{1}{n!} \\ \ln \lambda \\ n \end{matrix}$$

$$\text{Mult}(m_1, \dots, m_N | \mu, \bar{\mu}) = \frac{N!}{m_1! \dots m_N!} \exp \left[\sum m_k \log \mu_k \right]$$

FIXED

$$= " \exp \left[\sum_{k=1}^{N-1} m_k \log \mu_k + (N - m_1 - \dots) \log (1 - \mu_1 - \dots) \right]$$

$$\therefore h(x) = \frac{N!}{m_1! \dots m_{N-1}! (1 - \sum_{k=1}^{N-1} m_k)!}$$

$$u = \begin{pmatrix} m_1 \\ \vdots \\ m_{N-1} \end{pmatrix}$$

$$\eta_k = \log \frac{\mu_k}{1 - \sum_{i=1}^{N-1} \mu_i}$$

$$\left(\Leftrightarrow \mu_k = \frac{\exp \eta_k}{1 + \sum_{i=1}^{N-1} \exp \eta_i} \right)$$

- Formulae $\nabla_{\eta} \ln p(\eta) = \mathbb{E}[u(x)]$

- ML method

$$\begin{aligned}\eta_{ML} \text{ solves } 0 &= \nabla_{\eta} (p(x|\eta)) \\ &= \nabla_{\eta} \left[\prod_{n=1}^N h(x_n) \cdot g(\eta)^N \cdot \exp \left[\eta^T \cdot \sum_{n=1}^N u(x_n) \right] \right]\end{aligned}$$

$$\therefore \boxed{\nabla_{\eta} \ln p(\eta) \Big|_{\eta=\eta_{ML}} = \frac{1}{N} \sum_{n=1}^N u(x_n)}$$

sufficient statistics

$\sum u \rightarrow N x$
 $N \rightarrow n$

- Conjugate prior

$$p(\eta | X, v) = f(X, v) g(\eta)^v \exp \left[v \eta^T X \right]$$

$$p(\eta | X, \bar{X}, v) = \tilde{f}(X, \bar{X}, v) g(\eta)^{v+N} \exp \left[\eta^T (v \bar{X} + \sum u(x_n)) \right]$$

① Non-informative prior

- Constant prior? ... ① improper if unbound

② constant on which parametrization? ... take care!

- translation-invariant prior

... inv. under $X \rightarrow X + C$. has form $p(x|\mu) = f(x-\mu)$.

... $p(\mu) = \text{const.} \equiv \lim_{\sigma_0 \rightarrow 0} N(\mu | M_0, \sigma_0)$ ↑ location parameter

- scale-invariant prior

... inv. under $X \rightarrow k X$; form of $p(x|\sigma) = \frac{1}{\sigma} f\left(\frac{x}{\sigma}\right)$

... $p(\sigma) \propto 1/\sigma \equiv \lim_{\substack{a \rightarrow 0 \\ b \rightarrow 0}} \text{Gam}(\sigma | a, b)$ ↑ scale parameter

§ 2.5 Nonparametric methods

$$p(x) = \frac{K}{N V} \quad \begin{matrix} K \leftarrow \# \text{"near" events} \\ \# \text{events} \rightarrow \text{volume of "near"} \end{matrix}$$

- Histogram

... count events within bins. [careful on the bin-width! very bad for large dim. data]

- Kernel density estimators

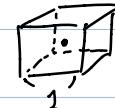
... count events that fall nearby.

location of n-th event

$$p(x) = \frac{1}{N} \cdot \frac{1}{h^D} \sum_{n=1}^N k\left(\frac{x - x_n}{h}\right)$$

OR $k(d) = \begin{cases} 1 : & i, d_i \leq \frac{1}{2} \\ 0 : & \text{otherwise} \end{cases}$

$$p(x) = \frac{1}{N (2\pi h^2)^{D/2}} \sum_{n=1}^N \exp\left[-\frac{\|x - x_n\|^2}{2 h^2}\right]$$



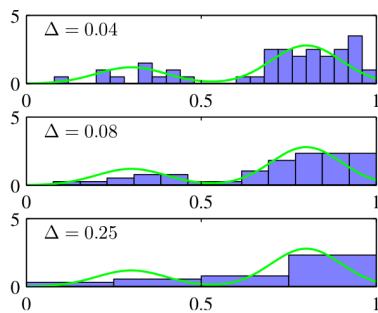
$$K(d) : K(d) \geq 0 \text{ and } \int K(d) d d = 1$$

- K-nearest neighbor method

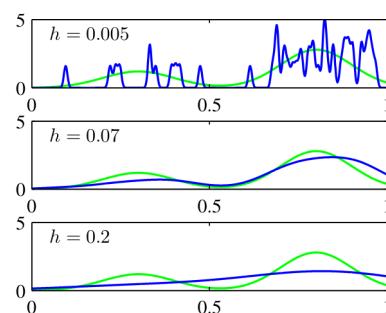
... calculate sphere volume V that includes K neighbors.

Needs entire data $\{x_n\}$ for evaluation !!

Histogram



Kernel density



K-nearest neighbor

