

# PRML 3. LINEAR MODELS FOR REGRESSION

Regression: Given  $\{(\mathbf{x}_n, t_n)\}_{n=1 \dots N}$ , estimate  $f$  for new  $\mathbf{x}$ .  
 $\dots f = y(\mathbf{x}) \text{ or } p(f|\mathbf{x})$

- Linear regression -- too simple

$$y(\mathbf{x}, \mathbf{w}) = w_0 + x_1 w_1 + \dots + x_D w_D$$

- Extended version

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{i=1}^{M-1} w_i \phi_i(\mathbf{x}) \quad \text{bias } \downarrow \quad \text{basis function } \mathbb{R}^D \rightarrow \mathbb{R}$$

$$\begin{aligned} &\equiv \sum_{i=0}^{M-1} w_i \phi_i(\mathbf{x}) \quad w. \mathbf{x} \in \mathbb{R}^D, \mathbf{w} \in \mathbb{R}^M \\ &= \mathbf{w}^\top \phi \quad \phi_0(\mathbf{x}) = 1. \end{aligned}$$

linear in  $\mathbf{w}$

Polynomial  $\phi_i(\mathbf{x}) = \mathbf{x}^i$  ... global :  $\mathbf{w}$  has long-range interaction

Gaussian  $\phi_i(\mathbf{x}) = \exp\left[-\frac{(\mathbf{x} - \mu_i)^2}{2s^2}\right]$

logistic sigmoid  $\phi_i(\mathbf{x}) = \sigma\left(\frac{\mathbf{x} - \mu_i}{s}\right)$   $\mu_i$  common [or  $\tanh \frac{\mathbf{x} - \mu_i}{s}$ ]

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

$\tanh a = 2\sigma(2a) - 1$

Fourier

:

Let's assume:  $t$  is given by  $t = y(x, w) + \epsilon \xleftarrow{\text{Gaussian noise}}$

$$\beta = 1/\sigma$$

i.e.  $P(t|x, w, \beta) = N(t|y(x, w), \beta^{-1})$  (precision)

$\Rightarrow$  Unimodality is assumed.

$$\mathbb{E}_t[t|x] = y(x, w)$$

Review: § 1.5

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

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

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

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

With data  $(t, \mathbf{x}) = (t_1, \mathbf{x}_1), \dots, (t_N, \mathbf{x}_N)$ ,

$$\ln P(t|\mathbf{x}, w, \beta) = \sum_{i=1}^N \ln N(t_i | y(\mathbf{x}_i, w), \beta^{-1})$$

$$\begin{aligned} &\downarrow \\ &\text{hereafter, omitted for simplicity} \quad = \frac{N}{2} \ln \frac{\beta}{2\pi} - \beta \left[ \frac{1}{2} \sum \left( t_n - w^\top \phi(\mathbf{x}_n) \right)^2 \right] \\ &=: E_0(w) \end{aligned}$$

Maximizing likelihood w.r.t.  $w$ ,

$$\nabla_w \ln P(t|w, \beta) = \beta \sum (t_n - w^\top \phi(\mathbf{x}_n)) \phi(\mathbf{x}_n) = 0$$

$$\therefore \sum_n t_n \phi_i(\mathbf{x}_n) = \sum_n \sum_k w_k \phi_k(\mathbf{x}_n) \phi_i(\mathbf{x}_n)$$

$$\boxed{A^\dagger := (A^\top A)^{-1} A^\top}$$

Moore-Penrose pseudo-inverse

$$\therefore w_{ML} = \Phi^\dagger t \quad \text{where } \Phi_{ij} = \phi_j(\mathbf{x}_i)$$

design matrix ( $N \times M$ )

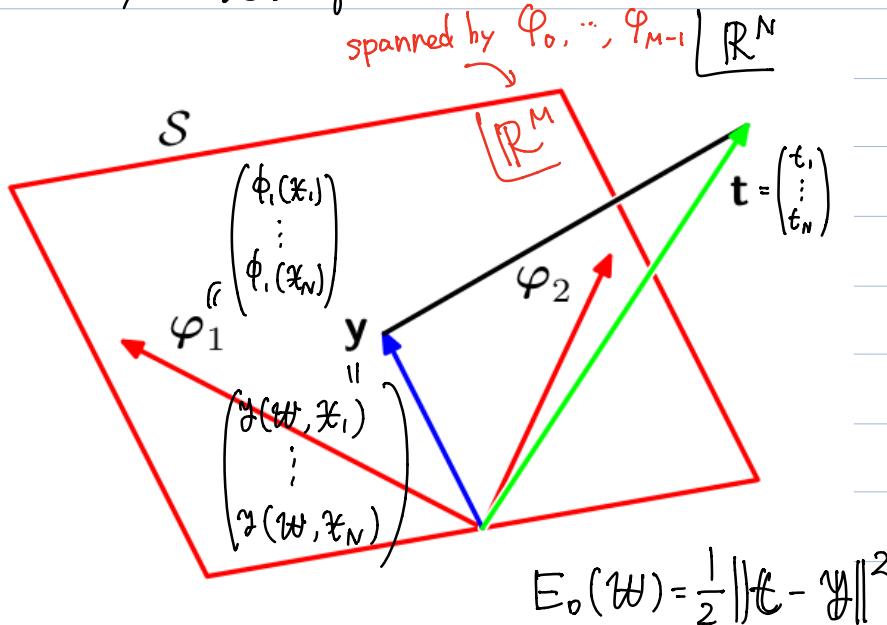
especially  $w_0^{ML} = \bar{t} - \sum_{m=1}^{M-1} \bar{w}_m \bar{\phi}_m$

avg. of  $t_n$       avg. of  $\phi_m(\mathbf{x}_i)$

As  $w_{ML}$  is indep. of  $\beta$ , we can calculate  $(w_{ML}, \beta_{ML})$  simply by  $\frac{\partial}{\partial \beta} \mathcal{L} = 0$ :

$$\beta_{ML}^{-1} = \frac{1}{N} \sum [t_i - w_{ML}^\top \phi(\mathbf{x}_i)]^2.$$

### § 3.1.2 Geometry of least squares



For  $M < N$ ,  $\{\varphi_0, \dots, \varphi_{M-1}\}$  spans a vector space  $S \subseteq \mathbb{R}^M \subseteq \mathbb{R}^N$

where  $y \in S$ ,  $t \in \mathbb{R}^N$ . so,  $y_{ML}$  is the projection of  $t$  onto  $S$ .  
this is  $\mathbb{P}^+ w_{ML}$ .

$$\begin{aligned} \text{PROOF } & \varphi_i \cdot (t - \mathbb{P} w_{ML}) \\ &= \varphi_i \cdot (t - \mathbb{P}(\mathbb{P}^\top \mathbb{P})^{-1} \mathbb{P}^\top t) \\ &= [\mathbb{P}^\top (t - \mathbb{P}(\mathbb{P}^\top \mathbb{P})^{-1} \mathbb{P}^\top t)]_i = 0, \end{aligned} \quad \therefore y_{ML} = \mathbb{P}^+ w_{ML}.$$

### § 3.1.3 Sequential learning

"stochastic gradient descent"

if  $E(w)$  is given by  $\sum_n E_n(w)$ ,

$$w^{(t+1)} = w^{(t)} - \eta \nabla_w E_n \Big|_{w=w^{(t)}} \quad \text{minimizes } E(w).$$

chosen so that  $w^{(t)}$  converges.

$$\left( \begin{array}{l} \text{for sum-of-square,} \\ w^{(t+1)} = w^{(t)} + \eta (t_n - w^{(t)\top} \phi_n) \phi_n \end{array} \right)$$

... least-mean-square algorithm.

### § 3.1.4 Regularized least squares

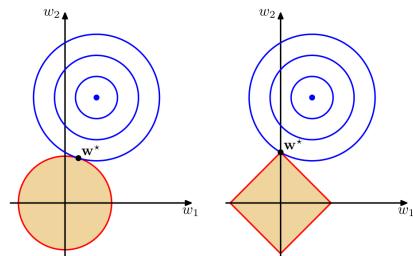
"weight decay"  $E_w(w) = \frac{1}{2} w^T w$ .

$$\text{e.g. } E(w) = \frac{1}{2} \sum_{n=1}^N [t_n - w^T \phi(x_n)]^2 + \frac{\lambda}{2} w^T w$$

$$\Rightarrow w = (\lambda I + \Phi^T \Phi)^{-1} \Phi^T t.$$

other regularizers:  $E_w(w) = \frac{\lambda}{2} \sum_{m=0}^{M-1} |w_m|^\beta$   
(or 1)

- Regularizers w.  $\beta \leq 1$   
drive some parameters to zero  
"sparse models"



### § 3.1.5 Multiple outputs $\cdots t \rightarrow \mathbf{t}$ .

- introduce extra basis functions for respective  $t_k$ .

$$\tilde{f}_k(x, w) = w^T \phi_k(x)$$

common for k

OR

- extend parameters  $w$  to  $\bar{w}$ .

$$\tilde{y}(x, \bar{w}) = \bar{w}^T \phi(x)$$

likelihood  $p(\mathbf{T} | \mathbf{X}, \bar{w}, \beta)$

$$\left[ p(t | x, \bar{w}, \beta) = N(t | \bar{w}^T \phi(x), \beta^{-1} I) \right.$$

returns  $\bar{w}_{\text{MC}} = (\Phi^T \Phi)^{-1} \Phi^T T$  as before.

$$\left. \quad \right]$$

## §3.2 Bias-Variance trade off

Review §1.5

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

$$L = [t - \tilde{y}(\mathbf{x})]^2 \text{ leads to } \tilde{y}(\mathbf{x}) = \mathbb{E}_t[t | \mathbf{x}]$$

$$\text{and } \mathbb{E}[L] = \underbrace{\int d\mathbf{x} [y(\mathbf{x}) - \tilde{y}(\mathbf{x})]^2 p(\mathbf{x})}_{①} + \underbrace{\int d\mathbf{x} dt [t - \tilde{y}(\mathbf{x})]^2 p(\mathbf{x}, t)}_{②}$$

Here we choose squared loss function.

$\left[ \begin{array}{l} \text{sum-of-square error function} \\ (\text{consequence of assuming a Gaussian noise}) \end{array} \right]$

② remains even if we obtained the best solution  $y = \tilde{y}$ .

Of course, based on limited data  $D$ ,

our estimated  $\hat{y} = \hat{y}(\mathbf{x}; D) \neq h(\mathbf{x})$ .

Note:  $\mathbb{E}_D[(y(\mathbf{x}; D) - h(\mathbf{x}))^2]$  depends on learning algorithms.

$$\Rightarrow \mathbb{E}_D[\textcircled{1}] = \mathbb{E}_D \left[ \left\{ \hat{y}(\mathbf{x}; D) - \tilde{y}(\mathbf{x}) \right\}^2 \right]$$

$$= \underbrace{\left\{ \tilde{y}(\mathbf{x}) - \mathbb{E}_D[\hat{y}(\mathbf{x}; D)] \right\}^2}_{\text{"bias"}^2} + \underbrace{\mathbb{E}_D \left[ \left\{ \hat{y}(\mathbf{x}; D) - \mathbb{E}_D[\hat{y}(\mathbf{x}; D)] \right\}^2 \right]}_{\text{"variance"}^2}$$

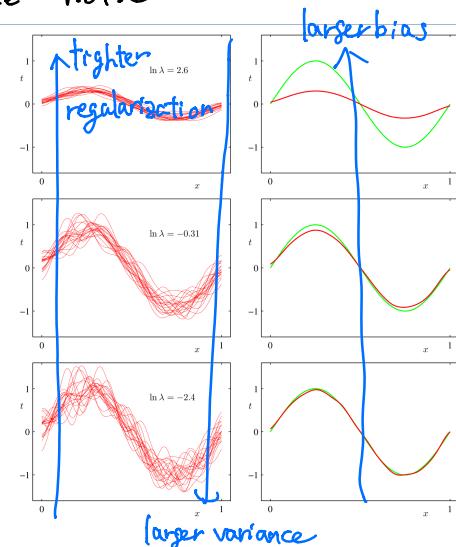
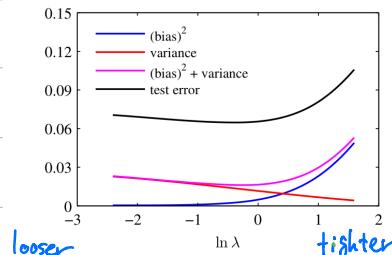
of our learning algorithm

∴ Expected loss  $\mathbb{E}_{D, \mathbf{x}}[L] = (\text{bias})^2 + \text{variance} + \text{noise}$

$$(\text{bias})^2 = \int \left\{ \mathbb{E}_D[\hat{y}(\mathbf{x}; D)] - \tilde{y}(\mathbf{x}) \right\}^2 p(\mathbf{x}) d\mathbf{x}$$

$$\text{variance} = \int \mathbb{E}_D \left[ \left\{ \hat{y}(\mathbf{x}; D) - \mathbb{E}_D[\hat{y}(\mathbf{x}; D)] \right\}^2 \right] p(\mathbf{x}) d\mathbf{x}$$

$$\text{noise} = \int [\tilde{y}(\mathbf{x}) - t]^2 p(\mathbf{x}, t) d\mathbf{x} dt$$



### §3.3 Bayesian linear regression

- ML method may lead to overfitting
- limit #basis ( $M$ ) ... less flexibility.
- use regulators ... which value of  $\lambda$ ?
- use Bayesian method ... needs proper prior

We use Gaussian likelihood ... conjugate prior is Gaussian.

prior:  $p(w) = N(w | M_0, S_0)$

posterior:  $p(w | t) = N(w | M_N, S_N)$

$$\begin{cases} M_N = S_N^{-1} M_0 + \beta \Phi^T t \\ S_N^{-1} = S_0^{-1} + \beta \Phi^T \Phi \end{cases}$$

• non-informative prior  $S_0 \rightarrow \infty \mathbb{I}$  gives  $M_N \rightarrow \Phi^T t$ ,  $S_N \rightarrow \beta \Phi^T \Phi$

A simple prior  $p(w | \alpha) = N(w | 0, \alpha^{-1} \mathbb{I}) \approx \exp\left[-\frac{\alpha}{2} w^T w\right]$

... equivalent to least-square error + squared regulator

$$\dots M_N = \beta S_N^{-1} \Phi^T t, \quad \ln p(w | t) = -\frac{\beta}{2} \sum_{n=1}^N [t_n - w^T \phi(x_n)]^2 - \frac{\alpha}{2} w^T w + \text{const.}$$

$$S_N = \alpha \mathbb{I} + \beta \Phi^T \Phi,$$

Another prior

$$p(w | \alpha) = \left[ \frac{q}{2 P(\gamma)} \left( \frac{\alpha}{2} \right)^{\frac{1}{\gamma}} \right]^M \exp\left(-\frac{\alpha}{2} \sum_{m=0}^{M-1} |w_m|^{\gamma}\right)$$

$$\Rightarrow \ln p(w | t) = -\frac{\beta}{2} \sum [t_n - w^T \phi(x_n)]^2 - \frac{\alpha}{2} \sum_{m=0}^{M-1} |w_m|^\gamma$$

- Remember that our goal is to predict  $t'$  from  $\mathbf{x}'$ .

$$P(t|\mathbf{x}, \mathbf{t}, \mathbf{X}, \alpha, \beta) = \int p(t|\mathbf{x}, w, \beta) p(w|t, \mathbf{X}, \alpha, \beta) dw \quad \dots \textcircled{1}$$

noise: assumed known prior  
 ①      ②  
 $\mathcal{N}(t|\mathbf{y}(\mathbf{x}, w), \beta^{-1})$        $\mathcal{N}(w|\mathbf{m}_N, \mathbf{S}_N)$

$$(2.115) \quad \begin{aligned} A &= \Phi^T \\ b &= 0 \end{aligned} \quad = \mathcal{N}(t | \mathbf{m}_N^T \Phi(\mathbf{x}), \beta^{-1} + \Phi^T(\mathbf{x}) \mathbf{S}_N \Phi(\mathbf{x}))$$

→ 0 for infinite data  
 noise from ①      uncertainty of  $w$

**WARNING** If we use localized basis functions,

e.g. gaussian

regions away from the locations will have too small  $\Phi^T \mathbf{S}_N \Phi$ , and the model becomes too confident on its prediction.

This can be patched by "Gaussian process".

= constructing the equivalent kernel directly w/o using basis functions.

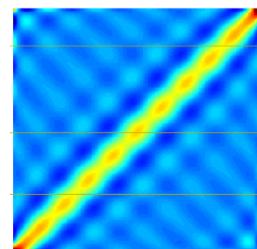
### § 3.3.3 Equivalent kernel

$$\begin{aligned} \mathbf{y}(\mathbf{x}, w) \Big|_{w=\mathbf{m}_N} &= \mathbf{m}_N^T \Phi(\mathbf{x}) \quad \xrightarrow{\text{mean of the predictive distribution}} P(t|\mathbf{x}, \mathbf{t}, \mathbf{X}, \alpha, \beta) \\ &= \beta \Phi(\mathbf{x}) \mathbf{S}_N \mathbf{T}^T \mathbf{t} = \sum_{n=1}^N [\underbrace{\beta \Phi(\mathbf{x})^T \mathbf{S}_N \Phi(\mathbf{x}_n)}_{=: k(\mathbf{x}, \mathbf{x}_n)}] t_n \\ &\quad \text{"smoother matrix" "equivalent kernel"} \end{aligned}$$

i.e. the mean of output  $t$  for an input  $\mathbf{x}$  is given by

$$\sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) t_n$$

↑      ↑      ↑  
 learned input/output  
 modelled by  $\Phi_i(\mathbf{x}_n)$ , etc...



Also:

$$\begin{aligned}\text{Cov}_{\mathbf{w}}[y(\mathbf{x}), y(\mathbf{x}')] &= \mathbb{E} \left[ (\mathbf{w}^\top \phi - \mathbf{m}_N^\top \phi)(\mathbf{w}^\top \phi' - \mathbf{m}_N^\top \phi') \right] \\ &= \phi(\mathbf{x})^\top \mathbb{E}[\mathbf{w} \mathbf{w}^\top] \phi(\mathbf{x}') + \phi(\mathbf{x})^\top \mathbf{m}_N \mathbf{m}_N^\top \phi(\mathbf{x}') \\ &\quad - \phi(\mathbf{x})^\top \left( \mathbb{E}[\mathbf{w}]^\top \mathbf{m}_N + \mathbf{m}_N^\top \mathbb{E}[\mathbf{w}] \right) \phi(\mathbf{x}') \\ \boxed{\mathbf{p}(\mathbf{w}) \equiv N(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N)} \quad &\quad = \phi(\mathbf{x})^\top \left( \mathbb{E}[\mathbf{w} \mathbf{w}^\top] - \mathbf{m}_N \mathbf{m}_N^\top \right) \phi(\mathbf{x}') \\ &= \phi(\mathbf{x})^\top \mathbb{V}[\mathbf{w}] \phi(\mathbf{x}') \\ &= \phi(\mathbf{x})^\top \mathbf{S}_N \phi(\mathbf{x}') \\ &= \beta^{-1} k(\mathbf{x}, \mathbf{x}')\end{aligned}$$

"predictive means of nearby points are highly correlated"  
distant

$$\bullet \sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) = 1 \quad \text{③ PROOF? or do we need } \begin{cases} \alpha = 0 \\ \phi_i(\mathbf{x}) : \text{L.indep.} \\ N \geq M \text{ etc.?} \end{cases}$$

$$\bullet k(\mathbf{x}, \mathbf{x}') = \psi^\top(\mathbf{x}) \psi(\mathbf{x}') \quad \text{with} \quad \psi(\mathbf{x}) = \sqrt{\beta} \mathbf{S}_N^{1/2} \phi(\mathbf{x}).$$

## \* § 3.4 Bayesian model comparison

Suppose we have  $L$  models  $\{M_i\}$  and data  $\mathcal{D}$  is generated from one of them.

"a model refers to a prob. dist. over the observed data"

... a model  $M_i$  generates a data set  $\mathcal{D}$

with a probability  $\int_{\text{params}} p(\mathcal{D}|M_i, \text{params}) p(\text{params}|M_i)$

prior probability distribution  $p(M_i)$

$\downarrow$  posterior  $p(M_i|\mathcal{D}) \propto p(\mathcal{D}|M_i) p(M_i)$

which model is  
more likely?

"model evidence"  
"marginal likelihood"

...  $\int_{\text{params}} p(\mathcal{D}|M_i, \text{params}) p(\text{params}) d(\text{params})$   
(usual) likelihood

$\hat{*}$  Bayes factor  $\frac{p(\mathcal{D}|M_i)}{p(\mathcal{D}|M_j)}$

$\downarrow$  predictive distribution

$p(t|\mathcal{X}, \mathcal{D}) = \sum_i p(t|\mathcal{X}, M_i, \mathcal{D}) p(M_i|\mathcal{D})$  "model mixture"

or we can choose a model by hand:  $p(t|\mathcal{X}, \mathcal{D}) := p(t|\mathcal{X}, \tilde{M}, \mathcal{D})$

"model selection"

$$\text{Here, } \mathbb{E}_{\mathcal{D}} \left[ \ln \frac{p(\mathcal{D}|M_{\text{true}})}{p(\mathcal{D}|M_i)} \right] = \int p(\mathcal{D}|M_{\text{true}}) \ln \frac{p(\mathcal{D}|M_{\text{true}})}{p(\mathcal{D}|M_i)} d\mathcal{D}$$

$$= KL(p_{\text{true}} || p_i) > 0,$$

which proves the Bayes factor chooses the true model on average as far as  $\{M_i\}$  includes the true model  $M_{\text{true}}$ .

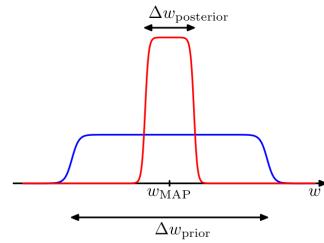
## • Interpretation of $P(D)$

Consider a model  $M$  that has one parameter  $w$ :

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

if we assume

the posterior  $P(w|D)$  is centered at  $w_{MAP}$  and flat  
and the prior  $P(w)$  is also flat,



$$P(D) \approx \Delta w_{post} P(D|M, w_{MAP}) \frac{1}{\Delta w_{prior}}$$

$$\therefore \ln P(D) \approx \ln P(D|M, w_{MAP}) - \ln \frac{\Delta w_{prior}}{\Delta w_{posterior}}$$

*fit result*

*→ 0 : larger penalty for tuned  $W$  (smaller  $\Delta w_{prior}$ )*  
*tuning penalty*

Note that the penalty could be larger for larger-parameter models.

Here we can see the prior must be proper (and thus informative).

Therefore, Bayesian approach also needs some validation data  
to validate the prior assumption

(instead of the model comparison.)

### §4.4.1 Bayesian Information Criterion

Using Laplace approx. to  $f = P(D|\Theta) P(\Theta)$ ,

$$\ln P(D) \approx \ln P(D|\Theta_{MAP}) + \underbrace{\ln P(\Theta_{MAP})}_{\text{Occam factor}} - \frac{1}{2} \ln |A| + \frac{M}{2} \ln 2\pi$$

with  $\Theta_{MAP}$  is the mode and

$$A = -\nabla \nabla \ln P(D|\Theta) P(\Theta)|_{MAP} = -\nabla \nabla P(\Theta_{MAP}|D).$$

If prior is broad and  $A$  is full rank,

$$\ln P(D) \approx \ln P(D|\Theta_{MAP}) - \frac{1}{2} M \ln N : \text{BIC.}$$

... BIC gives bad result when  $A$  is not full rank !!

### § 3.5 The Evidence Approximation

(empirical Bayes  
type-2 maximal likelihood  
generalized maximal likelihood)

Approximation for marginalizing  $\alpha$  and  $\beta$  as well as  $w$ .

hyperparameters      parameters

$$p(t | \mathbf{x}, \mathbf{t}, \mathbf{X}) = \int d\mathbf{w} d\alpha d\beta \underbrace{p(t | w, \beta)}_{N(t | y(w, \mathbf{x}), \beta^{-1})} \underbrace{p(w | \mathbf{t}, \alpha, \beta)}_{N(w | \mathbf{A}_N, \mathbf{S}_N)} p(\alpha, \beta | \mathbf{t})$$

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

likelihood  
(Evidence function)

- Analytic method

$$p(\alpha, \beta) \equiv \text{Gam}(\alpha | a, b) \text{Gam}(\beta | a', b') \quad \beta \mathbf{S}_N \mathbf{T} \mathbf{t} \quad [\alpha \mathbf{I} + \beta \mathbf{I}' \mathbf{I}]^{-1}$$

$$\Rightarrow p(t | \mathbf{t}) = \int d\mathbf{w} d\alpha d\beta \underbrace{N(t | y, \beta^{-1})}_{N(t | \hat{\alpha}, \hat{\beta}^{-1})} \underbrace{N(w | \mathbf{A}_N, \mathbf{S}_N)}_{N(w | \hat{w}, \hat{\mathbf{S}}_N)}$$

?

I guess

we need these approximations...

$$\approx \int d\mathbf{w} d\alpha d\beta \underbrace{N(t | y, \beta^{-1})}_{N(t | \hat{t}, \hat{\mathbf{S}}_t)} \underbrace{\frac{p(t | \hat{\alpha}, \hat{\beta})}{p(t)}}_{\text{irrelevant constant}} \\ \times \underbrace{N(w | \hat{w}, \hat{\mathbf{S}}_N \mathbf{T} \mathbf{t}, [\alpha \mathbf{I} + \beta \mathbf{I}' \mathbf{I}]^{-1})}_{\text{Student's } t} \underbrace{\text{Gam}(\alpha) \text{Gam}(\beta)}_{\text{Student's } t}$$

⊕ "Laplace approximation" for  $\int d\mathbf{w}$

(but this approx is known as poorer than the following ones.)

assuming  $p(\alpha, \beta)$  is less informative.

- Maximize  $p(\mathbf{t} | \alpha, \beta)$  [instead of  $p(\alpha, \beta | \mathbf{t})$ ]  $\rightarrow$  § 3.5.1

- expectation-maximization algorithm  $\rightarrow$  § 9.3.4

### § 3.5.1 - 2

$$P(\alpha, \beta | t) \propto P(t | \alpha, \beta) P(\alpha, \beta)$$

$$P(t | \alpha, \beta) = \int d\mathbf{w} P(t | \mathbf{w}, \beta) P(\mathbf{w} | \alpha)$$

$$= \int d\mathbf{w} \left[ \prod N(t_n | \mathbf{w}^T \phi(x_n), \beta^{-1}) \right] N(\mathbf{w} | 0, \alpha^{-1} \mathbb{I})$$

Gaussian noise assumption simple prior assumption

$$= \left( \frac{\beta}{2\pi} \right)^{N/2} \left( \frac{\alpha}{2\pi} \right)^{M/2} \int \exp[-E(\mathbf{w})] d\mathbf{w}$$

Here

$$E(\mathbf{w}) = \frac{\beta}{2} \|t - \Phi \mathbf{w}\|^2 + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w}$$

$$= E(\mathbf{m}_N) + \frac{1}{2} (\mathbf{w} - \mathbf{m}_N)^T A (\mathbf{w} - \mathbf{m}_N)$$

Showing that  $\mathbf{w}$  tends to be  $\mathbf{m}_N$  with an uncertainty  $S_N$ .

$$\dots E(\mathbf{m}_N) = \frac{\beta}{2} \|t\|^2 - \frac{1}{2} \mathbf{m}_N^T A \mathbf{m}_N$$

$$\text{with } A = \alpha \mathbb{I} + \beta \Phi^T \Phi = S_N^{-1}$$

$$\mathbf{m}_N = \beta A^{-1} \Phi^T t = \mathbf{m}_N \quad \left. \begin{array}{l} \text{or } P(\mathbf{w} | t) \\ \text{There } \alpha \text{ and } \beta \text{ were} \\ \text{"known" (assumed) parameters} \end{array} \right\}$$

$$A_{ij} = \frac{\partial}{\partial w_i} \frac{\partial}{\partial w_j} E(\mathbf{w}) \quad (\text{Hessian of } E)$$

$$\therefore \int \exp[-E(\mathbf{w})] d\mathbf{w} = e^{-E(\mathbf{m}_N)} (2\pi)^{M/2} |A|^{-1/2} \quad \text{and}$$

$$\ln P(t | \alpha, \beta) = \frac{N}{2} \ln \frac{\beta}{2\pi} + \frac{M}{2} \log \alpha - E(\mathbf{m}_N) - \frac{1}{2} \ln |A|.$$

We maximize  $P(t | \alpha, \beta)$  w.r.t.  $\alpha$  and  $\beta$ .

With the eigenvalues  $\{\lambda_i\}$  of  $\beta \Phi^T \Phi$  and  $\gamma = \sum_{i=1}^M \frac{\lambda_i}{\alpha + \lambda_i}$ ,

$$\frac{\partial P}{\partial \alpha} = 0 \Leftrightarrow \|\mathbf{m}_N\|^2 = \gamma/\alpha \quad (\lambda_i > 0)$$

$$\frac{\partial P}{\partial \beta} = 0 \Leftrightarrow \sum_{n=1}^N \left( t_n - \mathbf{m}_N^T \phi(x_n) \right)^2 = \frac{N - \gamma}{\beta}$$

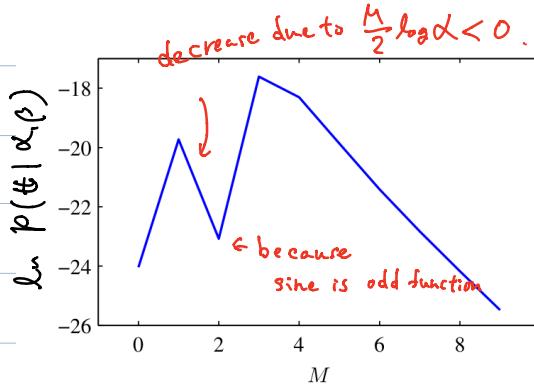
$A$  is positive definite  $\Leftrightarrow \forall x \neq 0, (Ax, x) > 0 \Leftrightarrow \text{all eigenvalues} > 0$ .  
(semi) (≥) (2)

$\therefore \beta \Phi^T \Phi$  is positive semi-definite.

Example:

to fit

with  $M$ -th polynomial



### Maximization

$$\mathbb{V}^T \mathbb{V} = \mathbb{V} \lambda \mathbb{V}^T \quad \text{eigenvalues } \lambda_i (= \lambda_i / \beta \text{ in text})$$

symmetric

$$[A = \mathbb{V} (\alpha \mathbb{I} + \beta \lambda) \mathbb{V}^T, A^{-1} = \mathbb{V} (\alpha \mathbb{I} + \beta \lambda)^{-1} \mathbb{V}^T]$$

$$\begin{aligned} E(\mathcal{M}_N) &= \frac{\beta}{2} \|\beta \mathbb{V} A^{-1} \mathbb{V}^T t - t\|^2 + \frac{\alpha \beta^2}{2} \|A^{-1} \mathbb{V}^T t\|^2 \\ &= t^T \left[ \frac{\beta}{2} (\beta \mathbb{V} A^{-1} \mathbb{V}^T - \mathbb{I}) (\beta \mathbb{V} A^{-1} \mathbb{V}^T - \mathbb{I})^T + \frac{\alpha \beta^2}{2} \mathbb{V} A^{-1} \mathbb{V}^T \right] t \\ &= \frac{\beta}{2} \mathbb{I} + \frac{\beta^3}{2} \mathbb{V} Z \lambda \mathbb{V}^T - \frac{\beta^2}{2} \mathbb{V} Z \mathbb{V}^T \mathbb{V}^T + \frac{\alpha \beta^2}{2} \mathbb{V} Z^2 \mathbb{V}^T \\ &\quad \xrightarrow{\text{using } \lambda_i = \frac{\beta}{\alpha + \beta \lambda_i}} \frac{\beta^2}{2} \frac{\beta M_i}{(\alpha + \beta M_i)^2} - \frac{\beta^2}{2} \frac{\beta}{\alpha + \beta M_i} + \frac{\beta^2}{2} \frac{\alpha}{(\alpha + \beta M_i)^2} = -\frac{\beta^2}{2} \mathbb{I} \\ &= \frac{\beta}{2} \mathbb{I} - \frac{\beta^2}{2} \mathbb{V} V (\alpha \mathbb{I} + \beta \lambda)^{-1} V^T \mathbb{V}^T \end{aligned}$$

$$\therefore \frac{\partial}{\partial \alpha} E(\mathcal{M}_N) = \frac{\beta^2}{2} t^T \mathbb{V} V (\alpha \mathbb{I} + \beta \lambda)^{-2} V^T \mathbb{V}^T t = \frac{1}{2} \|\mathcal{M}_N\|^2$$

$$\begin{aligned} \frac{\partial}{\partial \beta} E(\mathcal{M}_N) &= \frac{1}{2} \|t\|^2 + \frac{\beta^2}{2} t^T \mathbb{V} V Z^2 \lambda V^T \mathbb{V}^T t - \beta t^T \mathbb{V} V Z V^T \mathbb{V}^T t \\ &= \frac{1}{2} \|t\|^2 + \frac{1}{2} \|\mathbb{V} \mathcal{M}_N\|^2 - t^T \mathbb{V} \mathcal{M}_N \end{aligned}$$

$$\ln|A| = \sum \ln(\alpha + \beta \lambda_i) \quad \therefore \frac{\partial \ln|A|}{\partial \alpha} = \sum \frac{1}{\alpha + \lambda_i}, \frac{\partial \ln|A|}{\partial \beta} = \sum \frac{\lambda_i / \beta}{\alpha + \lambda_i}$$

$$\therefore \frac{\partial P}{\partial \alpha} = 0 \Leftrightarrow \alpha \|\mathcal{M}_N\|^2 = M - \sum \frac{\alpha}{\alpha + \lambda_i} = \sum \frac{\lambda_i}{\alpha + \lambda_i}$$

$$\frac{\partial P}{\partial \beta} = 0 \Leftrightarrow \beta \sum_{n=1}^N \left( t_n - \mathcal{M}_N^T \phi(x_n) \right)^2 = N - \sum \frac{\lambda_i}{\alpha + \lambda_i}$$

### § 3.5.3

maximum of the evidence function

$$P(t|\alpha, \beta) = \int d\mathbf{w} P(t|\mathbf{w}, \beta) p(\mathbf{w}|\alpha)$$

$$= \left( \frac{\beta}{2\pi} \right)^{N/2} \left( \frac{\alpha}{2\pi} \right)^{M/2} \int \exp[-E(\mathbf{w})] d\mathbf{w}$$

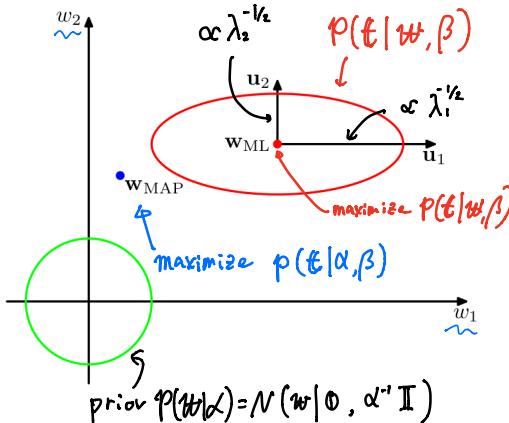
$$\begin{cases} \alpha = \frac{\gamma}{\| \mathbf{m}_N \|_2^2} \\ \beta^{-1} = \frac{1}{N-\gamma} \sum_{n=1}^N (t_n - \mathbf{m}_N^\top \phi(x_n))^2 \end{cases} \quad \text{where } \{\alpha + \lambda_i, \mathbf{u}_i\} \text{ is the eigensystem of the Hessian } A \text{ of } E(\mathbf{w}).$$

$$\gamma = \sum_{i=1}^M \frac{\lambda_i}{\alpha + \lambda_i} \quad \left( 0 < \frac{\lambda_i}{\alpha + \lambda_i} < 1 \right) \quad \text{effective number of well-determined parameters}$$

$$\mathbf{w}_{ML} = (\Phi^\top \Phi)^{-1} \Phi^\top t \Rightarrow \beta_{ML}^{-1} = \frac{1}{N} \sum (t_n - \mathbf{w}_{ML}^\top \phi(x_n))^2$$

$$\mathbf{w}_N = \mathbf{w}_{MAP} = (\frac{\alpha}{\beta} + \Phi^\top \Phi)^{-1} \Phi^\top t \Rightarrow \beta_{MAP}^{-1} = \frac{1}{N-\gamma} \sum (t_n - \mathbf{w}_{MAP}^\top \phi(x_n))^2$$

**Figure 3.15** Contours of the likelihood function (red) and the prior (green) in which the axes in parameter space have been rotated to align with the eigenvectors  $\mathbf{u}_i$  of the Hessian. For  $\alpha = 0$ , the mode of the posterior is given by the maximum likelihood solution  $\mathbf{w}_{ML} = \mathbf{m}_N$ . In the direction  $w_1$  the eigenvalue  $\lambda_1$ , defined by (3.87), is small compared with  $\alpha$  and so the quantity  $\lambda_1/(\lambda_1 + \alpha)$  is close to zero, and the corresponding MAP value of  $w_1$  is also close to zero. By contrast, in the direction  $w_2$  the eigenvalue  $\lambda_2$  is large compared with  $\alpha$  and so the quantity  $\lambda_2/(\lambda_2 + \alpha)$  is close to unity, and the MAP value of  $w_2$  is close to its maximum likelihood value.



$$E(\mathbf{w}) = E(\mathbf{m}_N) + \frac{1}{2} (\mathbf{w} - \mathbf{m}_N)^\top (\alpha \mathbf{I} + A) \mathbf{w}^\top (\mathbf{w} - \mathbf{m}_N)$$

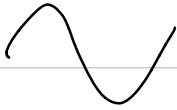
$$= E(\mathbf{m}_N) + \frac{1}{2} (\mathbf{w}' - \mathbf{m}'_N)^\top (\alpha \mathbf{I} + A) (\mathbf{w}' - \mathbf{m}'_N)$$

$$P(t|\mathbf{w}, \beta) = E(\mathbf{w}_N) + \frac{1}{2} (\mathbf{w} - \mathbf{w}_N)^\top [\beta \Phi^\top \Phi] (\mathbf{w} - \mathbf{w}_N)$$

$$= E(\mathbf{w}_N) + \frac{1}{2} (\mathbf{w}' - \mathbf{w}'_N)^\top A (\mathbf{w}' - \mathbf{w}'_N)$$

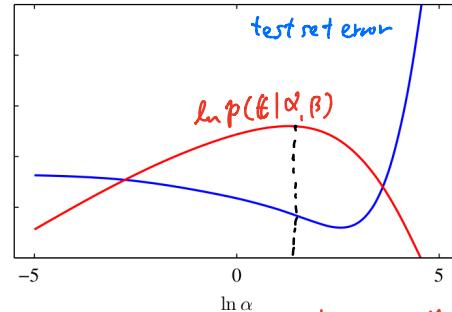
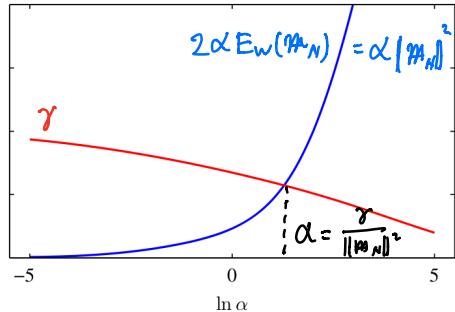
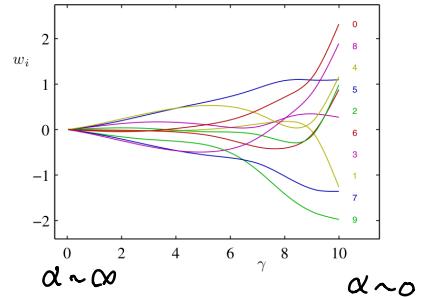
• Matrix  $A$  is diagonalizable  $\Leftrightarrow A \bar{A}^\top = \bar{A}^\top A$   
 $\Leftrightarrow$  a basis  $\{\mathbf{v}_i\}$  exists, where each  $\mathbf{v}_i$  is an eigenvector of  $A$   
Since Hessian is real & symmetric,  
its eigenvectors  $\{\mathbf{u}_i\}$  is a basis of  $\mathbb{R}^M$ .

Example: to fit



with 9-th polynomial

$(w_0, \dots, w_9)$  (with true  $\beta$ )



larger  $\alpha \Rightarrow$  tighter regularization  $\Rightarrow$  parameter is drawn less by data: smaller  $\gamma$

### § 3.6

Linear models = linear in  $W$

⊕ non-linear  $\phi(x)$  fixed before reading the data.

$\Downarrow$   
difficulty, especially in large dimension.

two properties to alleviate this problem:

① data tend to live in a lower-dimensional manifold

→ arrange  $\phi$  so that they are scattered only in the manifold

- support vector machine
- relevance vector machine

or adaptive  $\phi$

- neural network

②  $t$  tends to have strong dependency on specific direction of the manifold

... neural network can exploit this property.

# PRML 4. LINEAR MODELS FOR CLASSIFICATION

Classification: to separate the input space into  $K$  disjoint regions<sup>"decision regions"</sup>  
(possibly disconnected)  
i.e. to find  $(D-1)$ -dimensional hyperplane.

"decision boundary" "decision surface"

Three approaches

(a) calculate  $p(C_k | \mathbf{x})$ ,  $p(\mathbf{x})$ , etc. and then decide "generative model"  
... heavy computation

(b) calculate  $p(C_k | \mathbf{x})$  and then decide "discriminative model"

(c) find a function  $f: \mathbf{x} \mapsto k$  "discriminant function"  
~ no access to posterior probabilities.

- a target variable notation for probabilistic models

$\mathbf{f} = (P_1, \dots, P_K)^T$  where  $P_k = P(\mathbf{x} \in C_k)$   
or we can transform to  $\phi(\mathbf{x})$ .

$\Rightarrow$  "Generalized linear" model  $y(\mathbf{x}, \mathbf{w}) = f(\mathbf{x}^T \mathbf{w} + w_0)$   
activation function

## § 4.1.1 Discriminant functions for 2-class case ( $\mathbf{x} \mapsto 0 \text{ or } 1$ )

• linear discriminant  $y(\mathbf{x}) = w_0 + \mathbf{w}^T \mathbf{x}$  =  $\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}$ .  
bias weight vector

$\mathbf{x} \in C_1$  if  $y(\mathbf{x}) \geq 0$  ( $\mathbf{w}^T \mathbf{x} \geq -w_0$ )

$\mathbf{x} \in C_2$  if  $y(\mathbf{x}) < 0$  threshold

$\rightarrow$  decision plane  $P$  is  $D-1$ -dim. hyperplane

characterized by  $P \perp \mathbf{w}$ .  $\left( \text{and } d(0, P) = -\frac{w_0}{\|\mathbf{w}\|} \right)$

### § 4.1.2 Discriminant functions for multiple-class case

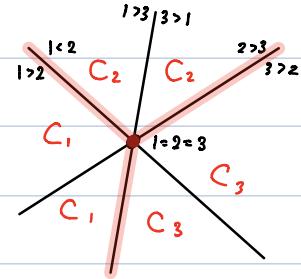
- one-versus-the-rest classifier ) do not work.
- one-versus-one classifier

$$y_k(x) = \mathbf{w}_k \cdot x + w_{k0}; \quad x \in C_k \text{ where } \forall j \neq k, y_j(x) > y_k(x)$$

decision boundary :  $y_k(x) = y_j(x)$

$$\dots (\mathbf{w}_k - \mathbf{w}_j)^T x + (w_{k0} - w_{j0}) = 0.$$

always singly connected & convex.



$$\begin{pmatrix} y_1 \\ \vdots \\ y_k \end{pmatrix} = \begin{pmatrix} \mathbf{w}_1 & \dots & \mathbf{w}_k \end{pmatrix}^T x + \begin{pmatrix} w_{01} \\ \vdots \\ w_{0k} \end{pmatrix} \quad \therefore y(x) = \tilde{\mathbf{W}}^T x + \tilde{w}_0 \\ = \tilde{\mathbf{W}}^T \tilde{x}$$

$$= \begin{pmatrix} w_{01} & w_{02} & \dots & w_{0k} \\ w_1 & w_2 & \dots & w_k \end{pmatrix} \begin{pmatrix} 1 \\ x \end{pmatrix} \quad \tilde{\mathbf{W}} = \begin{pmatrix} w_{01} & w_{02} & \dots & w_{0k} \\ w_1 & w_2 & \dots & w_k \end{pmatrix}$$

### § 4.1.3 Least square to determine $\tilde{\mathbf{W}}$ does not work

because least square method corresponds to MC under Gaussian assumption.

(and weak against outliers)

least square : for input  $\{x_n, t_n\}_{n=1, \dots, N}$   $\left[ \mathbb{X} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{pmatrix}, \mathbb{T} = \begin{pmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{pmatrix} \right]$

minimize  $E_0(\tilde{\mathbf{W}}) = \frac{1}{2} \text{Tr} [(\tilde{\mathbb{X}} \tilde{\mathbf{W}} - \mathbb{T})(\tilde{\mathbb{X}} \tilde{\mathbf{W}} - \mathbb{T})^T]$

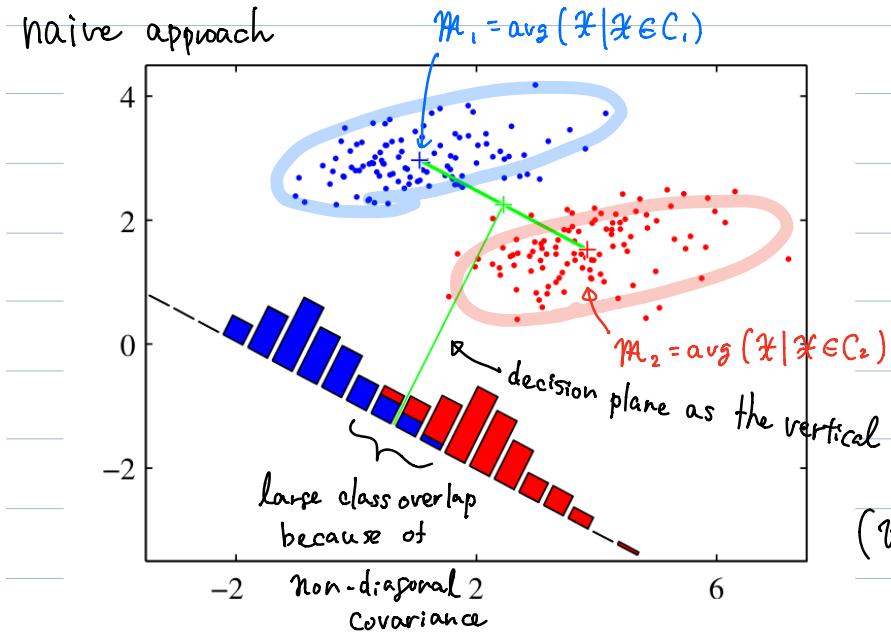
automatically satisfies  
 $\sum_{k=1}^K y_k(x) = 1$   
but not restricted to  $[0, 1]$

$$\dots \tilde{\mathbf{W}}_{\text{MC}} = (\tilde{\mathbb{X}}^T \tilde{\mathbb{X}})^{-1} \tilde{\mathbb{X}}^T \mathbb{T} = \tilde{\mathbb{X}}^+ \mathbb{T}. \quad y_{\text{MC}}(x) = (\tilde{\mathbb{X}}^+ \mathbb{T})^T \tilde{x}.$$

closed form, but very poor because  $t_n$  is far from Gaussian.

### § 4.1.4 Fisher's linear discriminant

naive approach



Fisher's method

$$\text{within-class variance } S_B^2(\mathbf{w}) = \sum_{n \in C_k} (\mathbf{w}^\top (\mathbf{x}_n - \mu_k))^2$$

⇒ Fisher criterion:

$$\text{maximize } J(\mathbf{w}) = \frac{(\mu_1 - \mu_2)^2}{S_1^2 + S_2^2} = \frac{[\mathbf{w}^\top (\mu_2 - \mu_1)]^2}{\sum_n [\mathbf{w}^\top (\mathbf{x}_n - \mu_k)]^2}$$

larger  $\mathbf{w} \cdot \Delta \mathbf{m}$   
...  $\mathbf{w} \parallel (\mu_1 - \mu_2)$

smaller  $\mathbf{w} \cdot (\mathbf{x}_n - \mu_k)$   
...  $\mathbf{w} \perp (\mu_1 - \mu_2)$

↑ for the class  $\mathbf{x}_n$  belongs to

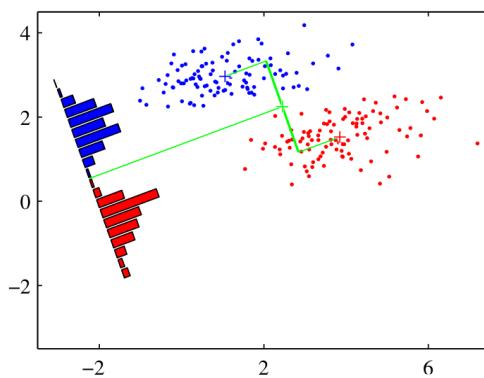
OR in matrix form.

$$J(\mathbf{w}) = \frac{\mathbf{w}^\top \mathbb{S}_B \mathbf{w}}{\mathbf{w}^\top \mathbb{S}_W \mathbf{w}}$$

$$\text{where } \mathbb{S}_B = (\mu_2 - \mu_1)(\mu_2 - \mu_1)^\top$$

$$\mathbb{S}_W = \sum_n (\mathbf{x}_n - \mu_k)(\mathbf{x}_n - \mu_k)^\top$$

for  $k$  to which  
 $\mathbf{x}_n$  belongs.



...  $J$  is maximized when  $(\mathbf{w}^\top \mathbb{S}_W \mathbf{w}) \mathbb{S}_B^{-1} \mathbf{w} = (\mathbf{w}^\top \mathbb{S}_B \mathbf{w}) \mathbb{S}_W^{-1} \mathbf{w}$  Fisher's linear discriminant

$\propto \mu_2 - \mu_1$

$\therefore \mathbf{w} \propto \mathbb{S}_W^{-1}(\mu_2 - \mu_1)$

§ 4.1.5 Fisher's solution is equivalent to least square  
under a certain target coding scheme

$$t = \begin{cases} N/N_1 & \text{for } C_1 \\ -N/N_2 & \text{for } C_2 \end{cases}$$

$$\Rightarrow E(\mathbf{w}) = \frac{1}{2} \left[ \sum_{n \in C_1} (\mathbf{w}^\top \mathbf{x}_n + w_0 - t_1)^2 + \sum_{n \in C_2} (\mathbf{w}^\top \mathbf{x}_n + w_0 - t_2)^2 \right]$$

$$\text{least square: } \sum_n (\mathbf{w}^\top \mathbf{x}_n) \mathbf{x}_n + w_0 \sum_n \mathbf{x}_n = N(\mu_1 - \mu_2)$$

$$w_0 = -\mathbf{w}^\top \mu \quad \text{where } \mu = \text{avg}(\mathbf{x}) = \frac{N_1}{N} \mu_1 + \frac{N_2}{N} \mu_2$$

$$\Rightarrow \left( \mathbb{S}_w + \frac{N_1 N_2}{N} \mathbb{S}_B \right) \mathbf{w} = N(\mu_1 - \mu_2)$$

$$\therefore \begin{cases} \mathbf{w} \propto \mathbb{S}_w^{-1} (\mu_2 - \mu_1) \\ w_0 = -\mathbf{w}^\top \mu \end{cases} \quad \begin{aligned} \mathbb{S}_w + \frac{N_1 N_2}{N} \mathbb{S}_B &= \sum \mathbf{x}_n \mathbf{x}_n^\top - N_1 \mu_1 \mu_1^\top - N_2 \mu_2 \mu_2^\top + \frac{N_1 N_2}{N} (\mu_1 - \mu_2)(\mu_1 - \mu_2)^\top \\ &= \sum \mathbf{x}_n \mathbf{x}_n^\top - N \mu \mu^\top. \end{aligned}$$

§ 4.1.6 Fisher's discriminant for ( $K > 2$ )-classes [with ( $D > K$ ) - input dimension]

Consider mapping  $\mathbf{x}_n$  to  $D'$  "features"  $\mathbf{y}_n = \mathbf{W}^\top \mathbf{x}_n \quad (\mathbf{y}_n \in \mathbb{R}^{D'}, \mathbf{x}_n \in \mathbb{R}^D)$   
no "bias" included

$$\mathbb{S}_w = \sum_{n=1}^{N_1} (\mathbf{y}_n - \mu_{C_1}) (\mathbf{y}_n - \mu_{C_1})^\top$$

$D' \times D'$

↑  
for  $C_k$  to which  $\mathbf{x}_n$  belongs

$$\mathbb{S}_B = \sum_{n=1}^N (\mu - \mu_{C_k}) (\mu - \mu_{C_k})^\top = \sum_{k=1}^K N_k (\mu_{C_k} - \mu) (\mu_{C_k} - \mu)^\top \quad \textcircled{1}$$

$D' \times D'$

$$\text{where } \mu_{C_k} = \text{avg}(\mathbf{y}_n | \mathbf{x}_n \in C_k), \mu = \text{avg}(\mathbf{y}_n) = \sum_{k=1}^K \frac{N_k}{N} \mu_{C_k}.$$

⇒ We are to find  $\mathbf{W}$  that maximizes  $J(\mathbf{W})$  ( $= \text{Tr}[\mathbb{S}_w^{-1} \mathbb{S}_B]$ , for example.)

\* Note that  $\text{rank}(\textcircled{1}) = 1$  and thus  $\text{rank}(\mathbb{S}_B) \leq K-1$ .

So we cannot find more than  $K-1$  features by this method.

...  $D'$  should be at most  $K$  (or  $K-1$ ).

## § 4.1.7 Perceptron

*bias  $\phi_0(\mathbf{x}) = 1$  is often included*

$$y(\mathbf{x}) = f(\mathbf{w}^T \phi(\mathbf{x})) \quad [2\text{-class}]$$

$$\text{where } f(a) = \begin{cases} +1 & (a \geq 0) \\ -1 & (a < 0) \end{cases} \quad (\text{also } t=1 \text{ for } C_1, t=-1 \text{ for } C_2)$$

Error function

"perceptron criterion"

$$E_p(\mathbf{w}) = - \sum_{n \in M} \mathbf{w}^T \phi_n t_n$$

*negative for  $n \in M$*

*positive for  $n \in N$  correctly classified*

$$\phi_n := \phi(\mathbf{x}_n)$$

$$M = \{n \in N \mid \text{misclassified}\}$$

"stochastic gradient descent"

if  $E(\mathbf{w})$  is given by  $\sum_n E_n(\mathbf{w})$ ,

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \gamma \nabla_{\mathbf{w}} E_n |_{\mathbf{w}=\mathbf{w}^{(t)}} \quad \text{minimizes } E(\mathbf{w}).$$

*chosen so that  $\mathbf{w}^{(t)}$  converges.*

$$\Rightarrow \mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \phi_n t_n \quad \text{for a misclassified } \mathbf{x}_n. \quad [\gamma \equiv 1 \text{ since } f(\gamma \mathbf{x}) = f(\mathbf{x})]$$

- perception converges  $\Leftrightarrow$  data is linearly-separable

Cons ) • We cannot determine the separability until it converges.

• Even if separable, the result depends on  $\mathbf{w}^{(0)}$  and the data order

• It won't converge for linearly-non-separable.

• No probabilistic result.

• Not ready to extend for multiple classes.

## §4.2 Probabilistic Generative Models

For 2-class case,

$$P(C_1 | \mathbf{x}) = \frac{P(\mathbf{x} | C_1) P(C_1)}{P(\mathbf{x})} = \frac{P(\mathbf{x} | C_1) P(C_1)}{P(\mathbf{x} | C_1) P(C_1) + P(\mathbf{x} | C_2) P(C_2)} = \frac{1}{1 + e^{-\alpha}} \stackrel{\text{logistic sigmoid}}{=} \sigma(\alpha)$$

where  $\alpha = \ln \frac{P(\mathbf{x} | C_1) P(C_1)}{P(\mathbf{x} | C_2) P(C_2)}$  ( $\sigma(\alpha) + \sigma(-\alpha) = 1$ )

For K-classes:  $P(C_k | \mathbf{x}) = \frac{\exp(\alpha_k)}{\sum \exp(\alpha_i)}$  where  $\alpha_k = \ln \left[ P(\mathbf{x} | C_k) P(C_k) \right]$

§4.2.1 Assume that  $\mathbf{x} \in C_k$  comes from  $N(\mathbf{x} | \mu_k, \Sigma)$ .

↑ assumed universal

$$P(\mathbf{x} | C_k) = \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \mu_k)^T \Sigma^{-1} (\mathbf{x} - \mu_k) \right]$$

For  $K=2$ ,

$$\begin{aligned} \alpha &= \ln \frac{P(C_1)}{P(C_2)} - \frac{1}{2} \left[ (\mathbf{x} - \mu_1)^T \Sigma^{-1} (\mathbf{x} - \mu_1) - (\mathbf{x} - \mu_2)^T \Sigma^{-1} (\mathbf{x} - \mu_2) \right] \\ P(C_1 | \mathbf{x}) &= \sigma(\mathbf{x}^T \mathbf{w} + w_0) = \mathbf{x}^T \underbrace{\Sigma^{-1}(\mu_1 - \mu_2)}_{=: \mathbf{w}} - \frac{1}{2} (\mu_1^T \Sigma^{-1} \mu_1 - \mu_2^T \Sigma^{-1} \mu_2) + \ln \frac{P(C_1)}{P(C_2)} =: w_0 \end{aligned}$$

For general K,

$$\begin{aligned} \alpha_k &= -\frac{1}{2} (\mathbf{x} - \mu_k)^T \Sigma^{-1} (\mathbf{x} - \mu_k) + \ln P(C_k) + \text{universal const.} \\ &\equiv \mathbf{x}^T \underbrace{\Sigma^{-1} \mu_k}_{\mathbf{w}_k} - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \ln P(C_k) \\ &=: \mathbf{x}^T \mathbf{w}_k + w_{0,k} \end{aligned}$$

$\therefore \Sigma$  is universal  $\Rightarrow$  boundary is linear in  $\mathbf{x}$  (generalized linear model)  
 non-universal  $\Rightarrow$  quadratic decision boundary

### § 4.2.2 ML solution \*Weak against outliers.

For  $K=2$ , we let  $t = \frac{1}{0}$  for  $\frac{C_1}{C_2}$  and prior  $p(C_1) = \frac{\pi}{1-\pi}$ .

$$\text{Likelihood } P(\mathbf{x}, \Sigma | \pi, \mu_1, \mu_2, \Sigma) \stackrel{\text{assumed universal}}{\sim} \prod_{n=1}^N \left[ \pi N(x_n | \mu_1, \Sigma) \right]^{t_n} \left[ (1-\pi) N(x_n | \mu_2, \Sigma) \right]^{1-t_n}$$

$\log P = N_1 \log \pi + N_2 \log (1-\pi) - \frac{N}{2} \log |\Sigma| - \sum_{n=1}^N \frac{1}{2} (x_n - \mu_{k_n})^\top \Sigma^{-1} (x_n - \mu_{k_n}) + \text{Const.}$

$$\Rightarrow \bar{\pi}_{\text{ML}} = \frac{N_1}{N_1 + N_2}, \quad (\mu_k)_{\text{ML}} = \arg \max_{\mu_k} P(x_n | x_n \in C_k)$$

$$\Sigma_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{k_n})(x_n - \mu_{k_n})^\top$$

$$\frac{\partial |A|}{\partial A} = |A| (A^\top)^{-1}$$

$$\frac{\partial A^{-1}}{\partial x} = -A^{-1} \frac{\partial A}{\partial x} A^{-1}$$

### § 4.2.3

For discrete input  $x \in \{0,1\}^D$ , general  $p(x|C_k)$  has  $2^D - 1$  parameters  $\theta_{k,d}$  for each class.

$$\text{If we simply assume } p(x|C_k) = \prod_{d=1}^D \mu_{k,d}^{x_d} (1-\mu_{k,d})^{1-x_d}$$

each  $x_d$  is independent with Bernoulli dist.  $\Rightarrow D$  parameter for each class

$$\begin{aligned} \text{Then } P(C_k|x) &= \frac{\exp(\alpha_k)}{\sum \exp(\alpha_j)} \text{ with } \alpha_k = \ln [P(x|C_k) P(C_k)] \\ &= \ln P(C_k) + \sum_{d=1}^D \left[ x_d \ln \frac{\mu_{k,d}}{1-\mu_{k,d}} + \ln(1-\mu_{k,d}) \right] \\ &\quad (\text{linear in } x_d !) \end{aligned}$$

### § 4.2.4

If class-conditional density  $p(x|C_k)$  is in exponential family with  $\mathcal{U}(x) = x$ ,  $\alpha_k$  is linear in  $x$ .

$$\begin{aligned} p(x|C_k) &= h(x) g(\theta_k) \exp(\theta_k^\top \mathcal{U}(x)) \\ \Rightarrow \alpha_k &= \underbrace{\theta_k^\top x}_{\text{Common and irrelevant.}} + \ln h(x) + \ln g(\theta_k) + \ln P(C_k) \end{aligned}$$

### § 4.3 Probabilistic Discriminative model

determine  $P(C_k | \mathbf{x})$  without  $P(\mathbf{x} | C_k)$  etc.

... directly by ML. "iterative reweighted least square"

Here we introduce fixed basis functions  $\phi(\mathbf{x})$

#### § 4.3.2

$$E[t | \phi] = p(t=1 | \phi) = y(\phi)$$

$$\text{logistic regression (for 2 classes)} \quad \begin{matrix} \text{logistic sigmoid} \\ \sigma = \frac{1}{1+e^{-\alpha}}, \quad \frac{d\sigma}{d\alpha} = \sigma(1-\sigma) \end{matrix}$$

$$p(C_1 | \phi) = y(\phi) = \sigma(\mathbf{w}^\top \phi) \quad \begin{matrix} \text{includes "bias"} \\ \phi_0(\mathbf{x}) = 1 \end{matrix}$$

... M parameters for  $\phi(\mathbf{x}) \in \mathbb{R}^M$

$$\left[ \begin{matrix} \text{Generative model needs} \\ M_k, \Sigma, p(C_i) = 2M + \frac{M(M+1)}{2} + 1 \text{ parameters} \end{matrix} \right]$$

$$P(\mathbf{t} | \mathbf{X}, \mathbf{w}) = \prod y_n^{t_n} (1-y_n)^{1-t_n}$$

$$E(\mathbf{w}) = -\ln P(\mathbf{t} | \mathbf{X}, \mathbf{w}) = -\sum [t_n \ln y_n + (1-t_n) \ln (1-y_n)]$$

$$\Rightarrow \nabla_{\mathbf{w}} E(\mathbf{w}) = \sum \frac{t_n - y_n}{y_n(y_n-1)} \nabla y_n = \sum (y_n - t_n) \phi(\mathbf{x}_n)$$

(→ we can employ some sequential algorithm)

\* This ML method can cause Overfitting for linearly-separable data.

ML solution  $\tilde{\mathbf{w}}$  satisfies  $\sum (\hat{y}_n - t_n) \phi(\mathbf{x}_n) = 0$ .

If the data is separable,

$t_n = 1 \Leftrightarrow \hat{y}_n > 0.5$ , i.e. the hyperplane  $\mathbf{w}^\top \phi = 0$  is the ML result.

$t_n = 0 \Leftrightarrow \hat{y}_n < 0.5$

Then ML prefers larger  $\|\mathbf{w}\|$  and  $\sigma(\mathbf{w}^\top \phi) \approx \text{ReLU}(\mathbf{w}^\top \phi)$ ,

which is very weak for new data.

⇒ We need prior or regularization.

### § 4.3.3

#### Newton-Raphson method

To find  $f(\mathbf{x}) = 0$  for  $f: \mathbb{R}^m \rightarrow \mathbb{R}$ ,  
 using  $f_k(\mathbf{x}) \approx f_k(\mathbf{x}_0) + \frac{\partial f_k}{\partial x_k}(x_k - x_{0k}) \approx 0$ ,  
 $x_k := x_k - (J^{-1})_{kk} f_k(\mathbf{x})$ , where  $J_{ij} = \frac{\partial f_i}{\partial x_j}(\mathbf{x})$ .  
 $\cdots \mathbf{x} := \mathbf{x} - J^{-1}f(\mathbf{x})$ .

→ Maximization of  $F(\mathbf{x})$ , where  $F: \mathbb{R}^m \rightarrow \mathbb{R}$ , is  $\nabla F = 0$ , thus

$$x_k := x_k - H_{kk}^{-1} \frac{\partial}{\partial x_k} F(\mathbf{x}) \text{ where } H_{ij} = \frac{\partial^2 F}{\partial x_i \partial x_j}$$

$$\cdots \mathbf{x} := \mathbf{x} - H^{-1} \nabla F.$$

$$\nabla_w E(w) = \sum (y_n - t_n) \phi_n, \quad H_{ij} = \sum_n y_n(1-y_n) \phi_i(\mathbf{x}_n) \phi_j(\mathbf{x}_n)$$

$$\text{or } H = \Phi^T R \Phi \text{ with } \Phi_{ij} = \phi_j(\mathbf{x}_i), \quad R = \text{diag}(y_n(1-y_n)).$$

<sup>positive definite because E is convex. ... unique minimum.</sup> <sup>Variance</sup>

$$\Rightarrow \text{Newton-Raphson} \quad \hat{w} := w - (\Phi^T R \Phi)^{-1} \Phi^T (y - t)$$

$$= (\Phi^T R \Phi)^{-1} \Phi^T R z$$

$$\xrightarrow{\text{iterative reweighting}} \text{where } z = \Phi \hat{w} - R^{-1}(y - t)$$

Mean and variance for logistic regression model  $P(C_1|\Phi) = \sigma(w^T \Phi)$

$$\mathbb{E}[t] = P(C_1|\Phi) + 0 \times P(C_2|\Phi) = y$$

$$\mathbb{V}[t] = \mathbb{E}[t^2] - \mathbb{E}[t]^2 = y(1-y)$$

We can use this  $z_n$  directly.

$$\text{Considering } Q_n = \hat{w}^T \phi_n \text{ and } E(a) = - \sum [t_n \log y_n + (1-t_n) \log(1-y_n)],$$

$$a_n := Q_n - \frac{y_n - t_n}{y_n(1-y_n)} = z_n.$$

$$\begin{cases} \nabla_a E = y - t \\ H_{nn} = y_n(1-y_n) \end{cases}$$

### §4.3.4 Logistic regression for K classes

$$P(C_k | \phi) = y_k(\phi) = \frac{\exp(\alpha_k)}{\sum \exp(\alpha_i)} \quad \text{with } \alpha_k = \mathbf{w}_k^T \phi$$

$$P(T | \mathbf{w}_1, \dots, \mathbf{w}_K) = \prod_{n=1}^N \prod_{k=1}^K t_{nk}^{y_k(\phi_n)} \quad \begin{matrix} \leftarrow t_n : 1\text{-of-}K \text{ coding scheme} \\ \leftarrow t_{nk} \\ \rightarrow y_{nk} \end{matrix}$$

$$E(T | \mathbf{w}_1, \dots, \mathbf{w}_K) = - \sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk} \quad \text{cross entropy error function}$$

$$\begin{aligned} \nabla_{\mathbf{w}_i} E &= - \sum_{n=1}^N \left( t_{ni} - \sum_{k=1}^K t_{nk} y_i \right) \phi(x_n) & \frac{\partial y_k}{\partial \alpha_i} &= (\delta_{ik} - y_i) y_k \\ &= - \sum_{n=1}^N \left( t_{ni} - y_i(\phi_n) \right) \phi_n & &= \delta_{ik} y_k - y_i y_k \end{aligned}$$

### §4.3.5 Other activation functions?

• noisy threshold model: assign  $t_n = 1$  iff  $\alpha_n = \mathbf{w}^T \phi_n > \theta$

and  $\theta$  is "noisy", i.e. given by some PDF  $g(\theta)$ .

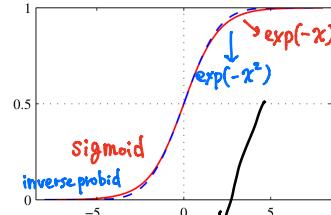
this results in  $P(t_n=1 | \alpha) = \int_{-\infty}^{\alpha} g(\theta) d\theta$ .

if we use  $g(\theta) = N(\theta | 0, 1)$ ,  $\xleftarrow{\text{equivalent to}} N(\theta | \mu, \sigma)$

CDF gives the inverse probit function

$$\begin{aligned} \Phi(a) &= \int_{-\infty}^a N(\theta | 0, 1) d\theta \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^a \exp(-\frac{1}{2}\theta^2) d\theta \end{aligned}$$

$$= \frac{1}{2} \left[ 1 + \operatorname{erf} \left( \frac{\mathbf{w}^T \phi}{\sqrt{2}} \right) \right]$$



inverse probit is weak against outliers.

$$\operatorname{erf}(a) = \frac{2}{\sqrt{\pi}} \int_0^a e^{-x^2} dx.$$

and we can determine  $\mathbf{w}$  by ML method.

• We can introduce a hyperparameter  $\epsilon$ : the probability that  $t$  has been mislabeled.

$$\text{Then } P(t|x) = (1-\epsilon) \delta(x) + \epsilon (1-\delta(x))$$

## ⑦ § 4.3.6 Canonical link functions

### § 4.2.4

If class-conditional density  $P(t|C_k)$  is in exponential family with  $U(x)=x$ ,  $\alpha_k$  is linear in  $x$ .

Consider "generalized linear model":

$$P(t|w, \phi) = \frac{1}{S} h\left(\frac{t}{S}\right) g(\eta) \exp\left(\frac{\eta t}{S}\right), \quad \text{where } \eta = \xi(w^\top \phi).$$

Then

$$\begin{aligned} y &\equiv \mathbb{E}[t|w, \phi] = -S \frac{d}{d\eta} \ln g(\eta) \Big|_{\eta=\xi(w^\top \phi)} & \Omega = \frac{\partial}{\partial \eta} \int P(t) dt = \int \left( \frac{g'(\eta)}{g(\eta)} + \frac{t}{S} \right) P(t) dt \\ &= : \psi^{-1}(\eta) \Big|_{\eta=\xi(w^\top \phi)} \quad : \therefore y = f(w^\top \phi) = (\psi^{-1} \circ \xi)(w^\top \phi). \end{aligned}$$

For GLM,

$$\begin{aligned} -\nabla_w \ln P(t|w, \phi) &= -\sum_{n=1}^N \left( \frac{t_n}{S} + \frac{d}{d\eta} \ln g(\eta_n) \right) \frac{\partial \eta}{\partial w} \\ &= \frac{1}{S} \sum_{n=1}^N (y_n - t_n) \nabla_w \xi(w^\top \phi_n) \end{aligned}$$

if we choose  $f$  so that  $\xi(a) = a$ ,

We can't choose  $f$ ?

$$\nabla_w \mathbb{E}(w) = \frac{1}{S} \sum_{n=1}^N (y_n - t_n) \phi_n.$$

I don't understand the importance / meaning of this section

If  $P(t|w, \phi) = \frac{t}{k} (1-k)^{1-t}$  with  $k = k(w^\top \phi)$ ,

$$\begin{aligned} P(t|w, \phi) &= (1-k) \exp\left[t \ln \frac{k}{1-k}\right] & \text{with } \begin{cases} \eta = \xi(w^\top \phi) = S \ln \frac{k}{1-k} \\ g(\eta) = S (1+e^\eta)^{-1} \end{cases} \\ &= (1+e^\eta)^{-1} \exp\left(\frac{\eta t}{S}\right) \end{aligned}$$

$$\begin{aligned} y &\equiv \mathbb{E}[t|w, \phi] = k(a) = (1+e^{-\eta/S})^{-1} & \dots \psi^{-1}(\eta) = (1+\exp(-\eta/S))^{-1} \\ &\dots f(a) = k(a). \end{aligned}$$



For Gaussian  $P(t|\mu) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(t-\mu)^2\right)$  with known  $\sigma$ ,

$$\begin{aligned} P(t|\mu) &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{t^2}{2\sigma^2}\right) \exp\left(-\frac{\mu^2}{2\sigma^2}\right) \exp\left(\frac{1}{\sigma^2}\mu t\right) & \begin{cases} \eta = \xi(a) = \frac{S}{\sigma^2} \mu(a) \\ g(\eta) = S \exp\left(-\frac{\sigma^2 \eta^2}{2S^2}\right) \exp\left(\frac{\eta t}{S}\right) \end{cases} \\ &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{S^2}{2\sigma^2} \frac{t^2}{S^2}\right) \exp\left(-\frac{\sigma^2 \eta^2}{2S^2}\right) \exp\left(\frac{\eta t}{S}\right) \end{aligned}$$

$$\mathbb{E}[t|\mu] = \mu = \frac{\sigma^2}{S} \eta, \quad \dots \psi^{-1}(\eta) = \frac{\sigma^2}{S} \eta, \quad f(a) = \mu(a).$$

$f(a) = \mu(a)$ . So what?

## § 4.4 Laplace approximation

... to give  $q(z) = \mathcal{N}(z | \mu, \sigma) \approx p(z)$  with  $\mu$  being the mode of  $p$ .

$$p(z) = \frac{f(z)}{Z}, \text{ where } Z = \int f(z) dz.$$

$$\ln f(z) \approx \ln f(z_0) - \frac{1}{2} \left[ -\frac{d^2}{dz^2} \ln f(z_0) \right] (z - z_0)^2, \text{ where } \begin{cases} f'(z_0) = 0 \\ f''(z_0) < 0. \end{cases}$$

$$\therefore p(z) = \frac{f(z)}{Z} \approx \sqrt{\frac{A}{2\pi}} \exp \left( -\frac{A}{2} (z - z_0)^2 \right),$$

$$A = -\frac{d^2}{dz^2} \ln f(z_0).$$

For multidimensional  $p(\mathbf{z})$ ,

$$p(\mathbf{z}) \approx \frac{|A|^{1/2}}{(2\pi)^{M/2}} \exp \left[ -\frac{1}{2} (\mathbf{z} - \mathbf{z}_0)^T A (\mathbf{z} - \mathbf{z}_0) \right]$$

$$\text{where } \nabla p = 0, \quad A = -\nabla \nabla \ln f(\mathbf{z}_0) \text{ [positive definite].}$$

$$\left[ \begin{array}{l} \text{also we can obtain} \\ \mathbf{z} \equiv \int d\mathbf{z} f(\mathbf{z}) \approx f(\mathbf{z}_0) \sqrt{\frac{(2\pi)^M}{|A|}} \left( = \frac{f(\mathbf{z}_0)}{\sqrt{A/2\pi}} \right) \end{array} \right].$$

## § 4.5 Bayesian logistic regression

$$P(w|t) \propto P(t|w) P(w)$$

$$\stackrel{\text{Bernoulli}}{\sim} \mathcal{N}(w|\mu_0, S_0)$$

$$\prod_{n=1}^N y_n^{t_n} (1-y_n)^{1-t_n} \text{ with } y_n = \sigma(w^\top \phi_n)$$

posterior

$$\Rightarrow \ln P(w|t) = -\frac{1}{2} (w - \mu_0)^\top S_0^{-1} (w - \mu_0) + \sum_{n=1}^N \left[ t_n \log y_n + (1-t_n) \log (1-y_n) \right] + C.$$

With Laplace approx,

$$P(w|t) \approx g(w|t) = \mathcal{N}(w|w_{MAP}, S_N)$$

where  $\mu_N = w_{MAP}$  is given by  $\nabla_w \ln P = 0$ , and

$$S_N^{-1} = -\nabla \nabla \ln P(w|t) = S_0^{-1} + \sum_{n=1}^N y_n (1-y_n) \phi_n \phi_n^\top.$$

Predictive distribution is

$$p(C_i | \phi, t) = \int d\omega p(C_i | \omega, \phi) P(\omega | t)$$

$$\approx \int da \sigma(a) \mathcal{N}(a | \mu_a, \sigma_a^2)$$

$$\text{with } \mu_a = w_{MAP}^\top \phi, \quad \sigma_a^2 = \phi^\top S_N \phi.$$

Also,  $\sigma(a) \approx \Phi(\sqrt{\frac{\pi}{8}} a)$  together with

$$\int_{-\infty}^{\infty} da \Phi(2a) \mathcal{N}(a | \mu, \sigma^2) = \Phi\left(\frac{\mu}{\sqrt{\sigma^2 + \lambda^{-2}}}\right)$$

$$\text{gives } p(C_i | \phi, t) \approx \sigma\left(\frac{\mu_a}{\sqrt{1 + \pi \sigma_a^2 / 8}}\right).$$

If we use the simplest decision boundary  $p(C_i | \phi, t) = 0.5$ ,

it is given by  $\mu_a = 0 \dots \phi \perp w_{MAP}$ .

$\Rightarrow$  equivalent to ML approach.