

# PRML 5. NEURAL NETWORKS

To adapt basis function to data.

- SVM : Select basis functions during training. ... Convex.
- RVM : ... sparser but non convex.
- NN : Fixed number, adapting during training. ... More compact. non convex.

## § 5.1

Feed-forward neural network

$$y_k(x, w) = F \left( \sum_{j=0}^M w_{kj}^{(2)} h(a_j) \right), \quad a_j = \sum_{i=0}^D w_{ji}^{(1)} x_i.$$

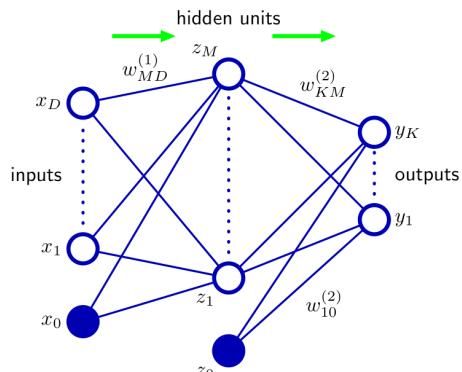
activation  
↓  
 $\phi$   
non-linear activation function

$f(x) = x$  for regression  
 $\sigma(x)$  for classification  
 (or softmax)

\* FNN with linear activation function  
 $\equiv$  linear transformation  
 (possibly with dimensionality reduction)

\* FNN may have skip-layers.

... sigmoidal hidden units can mimic them,  
 but including them can be better in practice.



- Two-layer FNN (except polynomial activation functions)

can uniformly approximate any continuous functions over a compact region.

- $M$  hidden units have  $2^M \times M!$ -wise symmetry. ( $+ \leftrightarrow -, i \leftrightarrow j$ )

→ at least  $2^M M!$  global minima.

## § 5.2 Natural choice of $F$ and $h$

$$y_{jk} = F(a_j^{(2)})$$

$$a_j^{(2)} = \sum_k w_{kj}^{(2)} z_j^{(1)}$$

$$z_j^{(1)} = h(a_j^{(1)})$$

$$a_j^{(1)} = \sum_i w_{ji}^{(1)} x_i$$

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

$$\text{ML approach: } \nabla_w E(w) = - \sum \nabla_w \ln P_n = 0$$

- Regression with Gaussian-noised data  $P(t|x, w) = N(t | y(x, w), \beta^{-1})$
- $\nabla_w \ln P = \beta(y - t) \nabla_w y$
- Classification  $P(t|x, w) = y^t (1-y)^{1-t}$
- $\nabla_w \ln P = \frac{y - t}{y(1-y)} \nabla_w y$
- Multi-class classification  $P(t_k|x, w) = \prod_k y_k^{t_k}$
- $\nabla_w \ln P = - \sum_k t_k \nabla_w \ln y_k$

$$\nabla_w y = \frac{\partial F}{\partial a^{(2)}} \nabla_w a^{(2)}$$

$$\sim \frac{\partial E}{\partial a_k}$$

$$\text{Regression: } F = \text{id.} \rightarrow -\nabla_w \ln P = \beta(y - t) \nabla_w a^{(2)}$$

$$\text{Classification: } F = \sigma \rightarrow F' = F(1-F) \rightarrow -\nabla_w \ln P = (y - t) \nabla_w a^{(2)}$$

$$\text{Multi-class: } F_k = \text{softmax} \rightarrow -\nabla_w \ln P = \sum_k (y_k - t_k) \nabla_w a_k^{(2)}$$

$$\dots \text{Choosing a suitable } F \text{ gives } \frac{\partial E}{\partial a_k} \propto y_k - t_k.$$

$$\begin{aligned} \frac{\partial}{\partial a_j} \ln F_k &= \delta_{jk} - F_j \\ -\nabla_w \ln P &= - \sum_k t_k \sum_j \left( \frac{\partial}{\partial a_j} \ln y_k \right) \nabla_w a_j^{(2)} \\ &= \sum_{j,k} t_k (y_j - \delta_{jk}) \nabla_w a_j^{(2)} \\ &= \sum_j \underbrace{[\sum_k t_k y_j - t_j]}_{=1} \nabla_w a_j^{(2)} \end{aligned}$$

## \* Gradient method

[https://www.slideshare.net/t\\_koshikawa/prml-5-pp227pp247](https://www.slideshare.net/t_koshikawa/prml-5-pp227pp247)

... but a few mistakes.

- global convergence ... converge to a solution regardless of the starting point

- convergence rate

$$p\text{-th order convergence} \Leftrightarrow \|\boldsymbol{x}_{k+1} - \boldsymbol{x}^*\| \leq C \|\boldsymbol{x}_k - \boldsymbol{x}^*\|^p. \quad (C < 1)$$

$$\text{super-linear convergence} \Leftrightarrow \|\boldsymbol{x}_{k+1} - \boldsymbol{x}^*\| \leq C_k \|\boldsymbol{x}_k - \boldsymbol{x}^*\|$$

with  $\lim_{k \rightarrow \infty} C_k = 0$ .

Linear << superlinear << quadratic.

- gradient descent

$$\boldsymbol{w} := \boldsymbol{w} - \mu \nabla E(\boldsymbol{w}) \quad \cdots \text{global, linear}$$

- Newton-Raphson

$$\boldsymbol{w} := \boldsymbol{w} - \alpha^{(t)} \underbrace{\boldsymbol{H}^{-1}}_{\text{tough, unstable}} \nabla E(\boldsymbol{w}) \quad \cdots \text{non-global, quadratic} \quad \text{!}$$

- Quasi-Newton

$$\boldsymbol{w}^{(t+1)} = \boldsymbol{w}^{(t)} - \alpha^{(t)} \underbrace{\boldsymbol{B}^{(t)}}_{\text{stable!}}^{-1} \nabla E(\boldsymbol{w}) \quad \text{with } \boldsymbol{B}^{(t)} \approx \boldsymbol{H},$$

$$\text{where } \boldsymbol{S} = \boldsymbol{w}^{(t)} - \boldsymbol{w}^{(t-1)}, \quad \boldsymbol{y} = \nabla E(\boldsymbol{w}^{(t)}) - \nabla E(\boldsymbol{w}^{(t-1)}), \quad \boldsymbol{\rho} = (\boldsymbol{y}^T \boldsymbol{S})^{-1},$$

$$\boldsymbol{B} := \boldsymbol{B} - \frac{\boldsymbol{B} \boldsymbol{S} \boldsymbol{S}^T \boldsymbol{B}}{\boldsymbol{S}^T \boldsymbol{B} \boldsymbol{S}} + \rho \boldsymbol{y} \boldsymbol{y}^T$$

$$\left[ \text{or } \boldsymbol{B}^{-1} := (\boldsymbol{I} - \rho \boldsymbol{S} \boldsymbol{y}^T) \boldsymbol{B}^{-1} (\boldsymbol{I} - \rho \boldsymbol{y} \boldsymbol{S}^T) + \rho \boldsymbol{S} \boldsymbol{S}^T \right] \quad (\text{BFGS})$$

- Conjugate descent.

### § 5.3 Back propagation to evaluate $\nabla E(w)$ .

Define  $\delta_j^{(m)} \equiv \frac{\partial E_n}{\partial a_j^{(m)}}$  where  $E(w) = \sum_n E_n(w)$ .

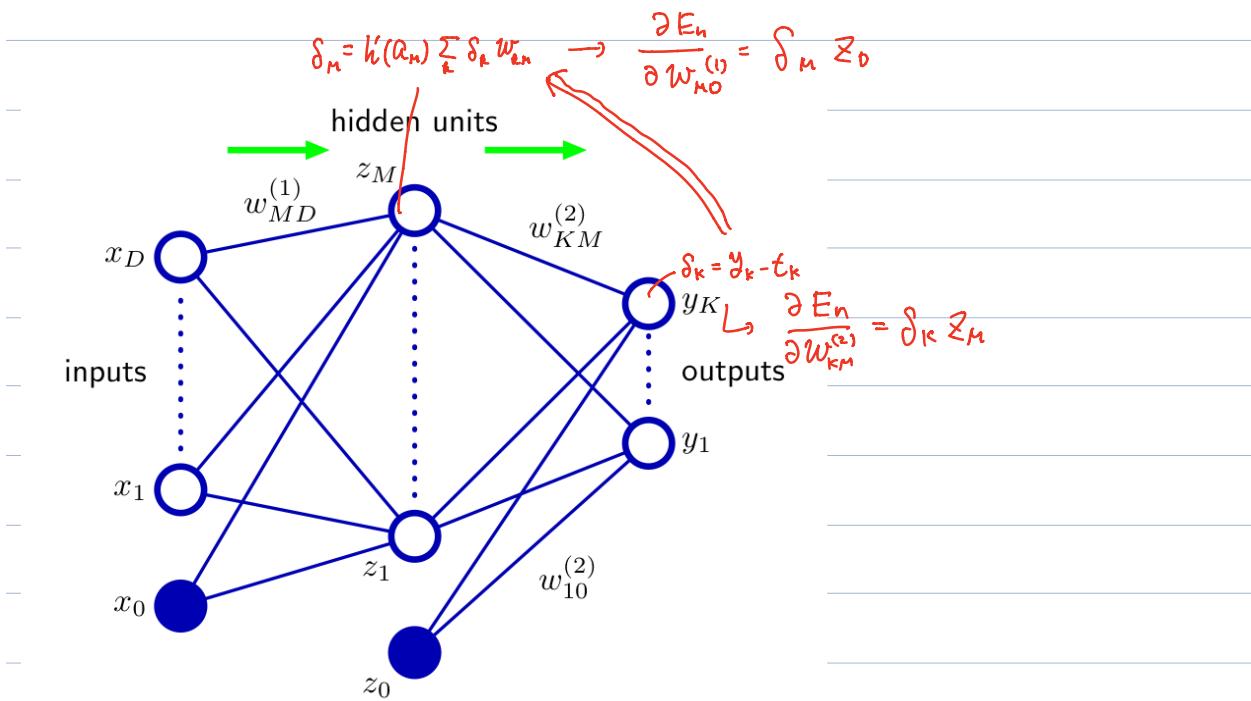
$$\frac{\partial E_n}{\partial w_{ji}^{(m)}} = \frac{\partial E_n}{\partial a_j^{(m)}} \frac{\partial a_j^{(m)}}{\partial w_{ji}^{(m)}} = \delta_j^{(m)} z_i^{(m-1)}$$

$$\begin{aligned}y_k &= f(a_k^{(s)}) \\a_k^{(s)} &= \sum w_{kj}^{(s)} z_j^{(t)} \\z_j^{(t)} &= h(a_j^{(t)}) \\a_j^{(t)} &= \sum_{i=0}^n w_{ji}^{(t)} x_i.\end{aligned}$$

For output units,  $\delta_j^{(out)} \equiv \frac{\partial E_n}{\partial a_j^{(out)}} = y_j - t_j$

while for hidden units,  $\delta_j^{(m)} \equiv \frac{\partial E_n}{\partial a_j^{(m)}} = \sum_k \frac{\partial E_n}{\partial a_k^{(m+1)}} \frac{\partial a_k^{(m+1)}}{\partial z_j^{(m)}} \frac{\partial z_j^{(m)}}{\partial a_j^{(m)}}$

$$= \sum_k \delta_k^{(m+1)} w_{kj}^{(m+1)} h'(a_j^{(m)}).$$



\* Numerical differentiation of  $\frac{\partial E}{\partial w_{ji}}$  is slower than Back propagation but useful to debug.

## § 5.4 Hessian

$$H_{(j;i)(l;k)} = \frac{\partial^2 E}{\partial w_{ji} \partial w_{lk}}$$

### § S.4.1 Diagonal approximation

Useful to calculate  $H^{-1}$ , but not good approximation.

$$\frac{\partial^2 E_n}{\partial W_{ji}^2} = \frac{\partial^2 E_n}{\partial a_j^2} Z_i^2 ;$$

$$\frac{\partial^2 E_n}{\partial a_j^2} = \frac{\partial}{\partial a_j} \left( \sum_k \frac{\partial E_n}{\partial a_k} \tilde{w}_{kj} h(a_j) \right)$$

$\uparrow$  CONNECT to NEXT LAYER  
NEXT LAYER

$\uparrow$  m-th layer

$$= \sum_k h''(\alpha_j) W_{kj} \frac{\partial E_n}{\partial \alpha_k} + h'(\alpha_j)^2 \sum_{k,k'} W_{kj} W_{k'j} \frac{\partial^2 E_n}{\partial \alpha_k \partial \alpha_{k'}} \\ \uparrow m+1-th \quad \quad \quad m+1-th \rightarrow \quad \quad \quad \} \text{neglecting off-diagonal}$$

### § 5.4.2 Outer product approximation (Levenberg - Marquardt approximation)

Valid only for appropriately-trained network.

$$\text{For } E_n = \frac{1}{2} (y_n - t_n)^2,$$

$$H = \nabla_w \nabla_w E = \sum_{n=1}^N \left[ (\nabla y_n)(\nabla y_n)^T + (y_n - t_n) \nabla \nabla y_n \right].$$

$$\approx \sum \|b_n\| b_n^* \quad \text{w.} \quad b_n = \nabla Q_n = \nabla Y_n.$$

For cross-entropy error  $E_n = -t_n \log y_n$  with  $\sigma(a)$ ,

$$H \approx \sum_{n=1}^N y_n (1-y_n) \|b_n\| b_n^T.$$

### § 5.4.3

$$H_N = \sum_{n=1}^N b_n b_n^\top$$

$$\Rightarrow H_N^{-1} = H_{N-1}^{-1} - \frac{H_{N-1}^{-1} b_N b_N^\top H_{N-1}^{-1}}{1 + b_N^\top H_{N-1}^{-1} b_N}$$

( $H_0$  can be set  $H_0 = \alpha I$  with small  $\alpha$ )

### § 5.4.4 Numerical differentiation

for  $f(x)$ ,

$$\frac{\partial f}{\partial x} = \frac{f(x+\varepsilon) - f(x-\varepsilon)}{2\varepsilon} + O(\varepsilon^2)$$

similarly, for  $f(x, y)$ ,

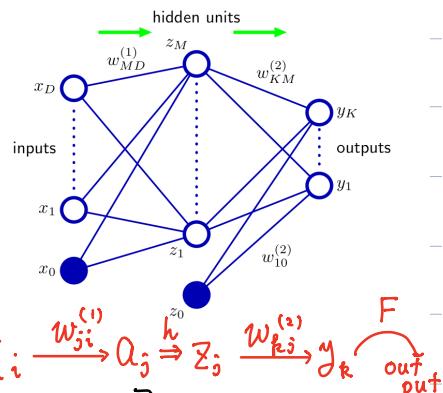
$$\begin{aligned} \frac{\partial^2 f}{\partial x \partial y} &= \frac{1}{4\varepsilon^2} \left[ f(x+\varepsilon, y+\varepsilon) - f(x+\varepsilon, y-\varepsilon) \right. \\ &\quad \left. - f(x-\varepsilon, y+\varepsilon) + f(x-\varepsilon, y-\varepsilon) \right] + O(\varepsilon^2) \\ &= \frac{1}{2\varepsilon} \left[ \frac{\partial f}{\partial x}(x, y+\varepsilon) - \frac{\partial f}{\partial x}(x, y-\varepsilon) \right] + O(\varepsilon^2) \end{aligned}$$

### § 5.4.5

$$\frac{\partial^2 E_n}{\partial w_{kj}^{(2)} \partial w_{k'j'}^{(2)}} = z_j z_{j'} M_{kk'}$$

$$\begin{aligned} \frac{\partial^2 E_n}{\partial w_{ji}^{(1)} \partial w_{j'i'}^{(1)}} &= x_i x_{i'} h''(a_{j'}) I_{jj'} \sum_k w_{kj'}^{(2)} \delta_k \\ &+ x_i x_{i'} h'(a_{j'}) h'(a_j) \sum_k \sum_{k'} w_{k'j'}^{(2)} w_{kj}^{(2)} M_{kk'} \end{aligned}$$

$$\frac{\partial^2 E_n}{\partial w_{ji}^{(1)} \partial w_{k'j'}^{(2)}} = \chi_i h'(a_j) \left[ \delta_k \mathbb{I}_{jj'} + z_{j'} \sum_{k'} w_{k'j'}^{(2)} M_{kk'} \right]$$



where  $\delta_k = \frac{\partial E_n}{\partial a_k}, \quad M_{kk'} \equiv \frac{\partial^2 E_n}{\partial a_k \partial a_{k'}}$

### § 5.4.6 Fast calculation of $HU$

Denoting  $\nabla^T \nabla_w = R$ , ← Because  $\nabla: \mathbb{R} \rightarrow \mathbb{R}^W$ ,  $R: \mathbb{R} \rightarrow \mathbb{R}$

$\overset{\text{D}}{\mathbb{R}^W}$

$$\nabla^T H = \nabla^T \nabla (\nabla E) = R(\nabla E). \quad \text{Note that } R(W) = U.$$

E.g. for 2-layer network with  $F = id.$ ,

$\nabla E$  is given by backpropagation as

$$\begin{cases} y_k = \sum w_{kj} z_j \\ z_j = h(a_j) \\ a_j = \sum w_{ji} x_i. \end{cases}$$

$$\frac{\partial E}{\partial w_{kj}} = \delta_k z_j = z_j (y_k - t_k),$$

$$\frac{\partial E}{\partial w_{ji}} = \delta_j x_i = x_i h'(a_j) \sum_k \delta_k w_{kj}.$$

$$\text{While, } R(y_k) = \sum_j w_{kj} R(z_j) + \sum_i v_{ki} z_j. \quad R(z_j) = h'(a_j) R(a_j)$$

$$R(\delta_k) = R(y_k)$$

$$R(a_j) = \sum_i v_{ji} x_i,$$

$$R(\delta_j) = \sum_k \left[ h''(a_j) R(a_j) w_{kj} \delta_k + h'(a_j) v_{kj} \delta_k + h'(a_j) w_{kj} R(\delta_k) \right]$$

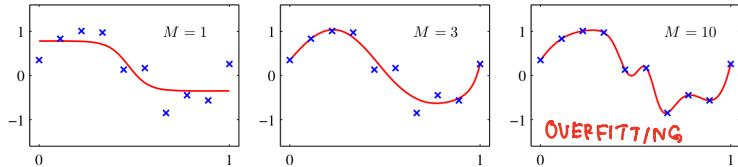
Thus  $\nabla^T H = R(\nabla E)$  is given by

$$R\left(\frac{\partial E}{\partial w_{kj}}\right) = R(\delta_k) z_j + \delta_k R(z_j)$$

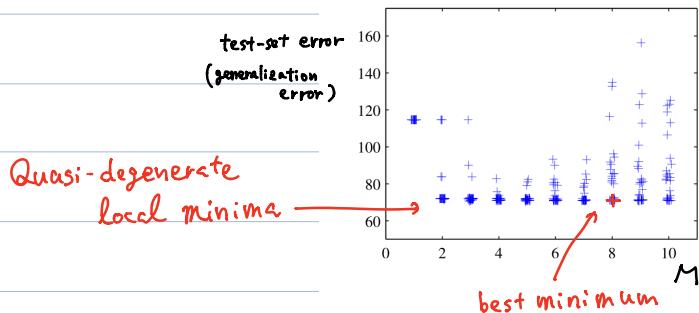
$$R\left(\frac{\partial E}{\partial w_{ji}}\right) = x_i R(\delta_j).$$

## § 5.5 Regularization

... How to determine  $\#\text{(hidden units)} \equiv M$ ?



"M giving the best test-set error" is not the optimal M.



### § 5.5.1

Simple weight-decay  $E_0 = \frac{\lambda}{2} \|w\|^2$  does not work  
because it does not respect the symmetry

$$\begin{cases} \mathbf{x} \mapsto A\mathbf{x} + \mathbf{b} \\ w_{ji} \mapsto w_{ja} A^{-1}_{ai} \\ w_{j0} \mapsto w_{j0} - w_{ja} A^{-1}_{ai} b_i \end{cases}$$

$$\text{OR} \quad \begin{cases} \mathbf{y} \mapsto A\mathbf{y} + \mathbf{b} \\ w_{kj} \mapsto A_{ka} w_{aj} \\ w_{k0} \mapsto A_{ka} w_{a0} + b_k. \end{cases}$$

The text does not discuss this general linear transformation but rather just a scale transformation.  
I suppose it is because we cannot build a generalizer that respects this symmetry.

We can build a generalizer that respects  $A = \alpha I$ , a subset of the above symmetry:

$$E_D = \frac{\alpha_1}{2} \sum_{w \in W^{(1)}} w^2 + \frac{\alpha_2}{2} \sum_{w \in W^{(2)}} w^2$$

This corresponds to the prior without bias weights.

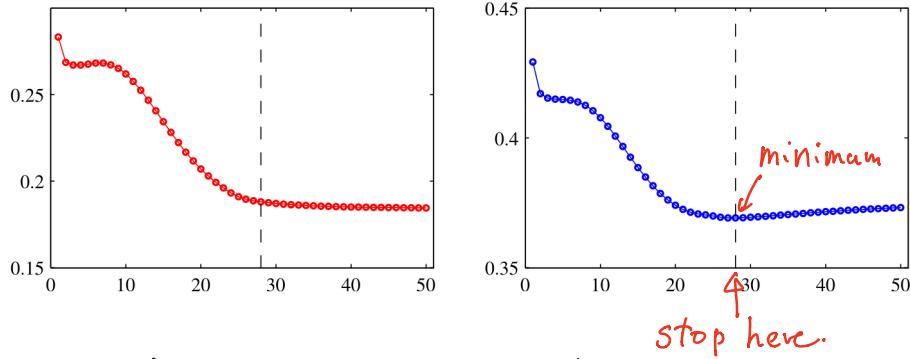
$$P(w | \alpha_1, \alpha_2) \propto \exp \left[ -\frac{\alpha_1}{2} \sum_{w \in W^{(1)}} w^2 - \frac{\alpha_2}{2} \sum_{w \in W^{(2)}} w^2 \right].$$

This is improper due to the absence of the biases,

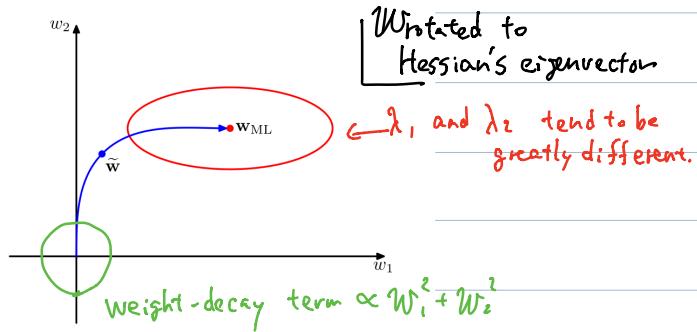
So it is common to include separate priors for biases.

Then no longer b-invariant.

## § 5.5.2 Early stopping



For a quadratic error function,



## § 5.5.3

Output should be stable against certain changes of input.

→ We want to tell the invariances to NN

(Otherwise much data are required for NN to learn the invariances by itself.)

- Multiplicate the data with the "invariant" changes.
- Add penalty to changes of output for transformed input: "tangent propagation".
- Preprocessing
- Invariant properties to be built into the structure of NN.

### § 5.5.4 Tangent propagation ... penalty for invariance.

Continuous transformation with L-parameter:

$$x \mapsto S(x, \xi) \stackrel{\text{infinitesimal}}{\simeq} x + \xi_L T_L(x); \quad T_L(x) = \frac{\delta S}{\delta \xi_L}(x, 0)$$

Considering  $L=1$  for simplicity, the output will be modified by

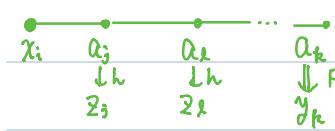
$$\frac{\delta y}{\delta \xi} \Big|_{\xi=0} = \sum_i \frac{\partial y}{\partial x_i} \frac{\partial x_i}{\partial \xi} \Big|_{\xi=0} = J T.$$

and we can penalize this:

$$E_{\text{transf.}} = \lambda \sum_{n=1}^N \Omega_n, \quad \text{where } \Omega_n = \frac{1}{2} \| J(x_n) T(x_n) \|^2.$$

Jacobian  $J = \frac{\partial y}{\partial x}$  can be calculated by backprop:

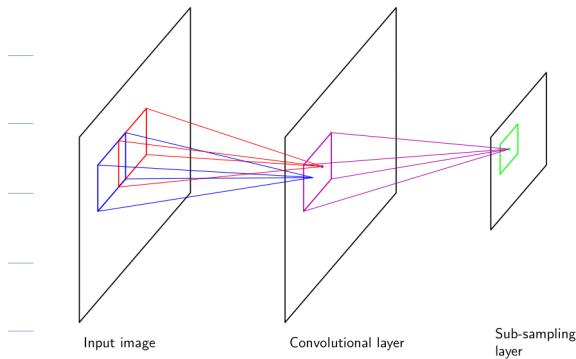
$$\begin{aligned} J_{ki} &= \frac{\partial y_k}{\partial x_i} = \sum_j \frac{\partial y_k}{\partial a_j} \frac{\partial a_j}{\partial x_i} & \frac{\partial y_k}{\partial a_j} &= \sum_l \frac{\partial y_k}{\partial a_l} \frac{\partial a_l}{\partial a_j} \\ &= \sum_j w_{ji} \frac{\partial y_k}{\partial a_j} & &= h'(a_j) \sum_l w_{lj} \frac{\partial y_k}{\partial a_l} \end{aligned}$$



and  $T(x_n)$  is calculated by

$$T(x) \simeq \frac{1}{2\varepsilon} [S(x, \varepsilon) - S(x, -\varepsilon)].$$

## § 5.5.6 Convolutional Neural Network



Invariant properties to be built into the structure of NN.

→ We need to know the "symmetries" (e.g. translation invariance) beforehand to use this technique.

a convolutional layer = feature maps

a feature map has units :

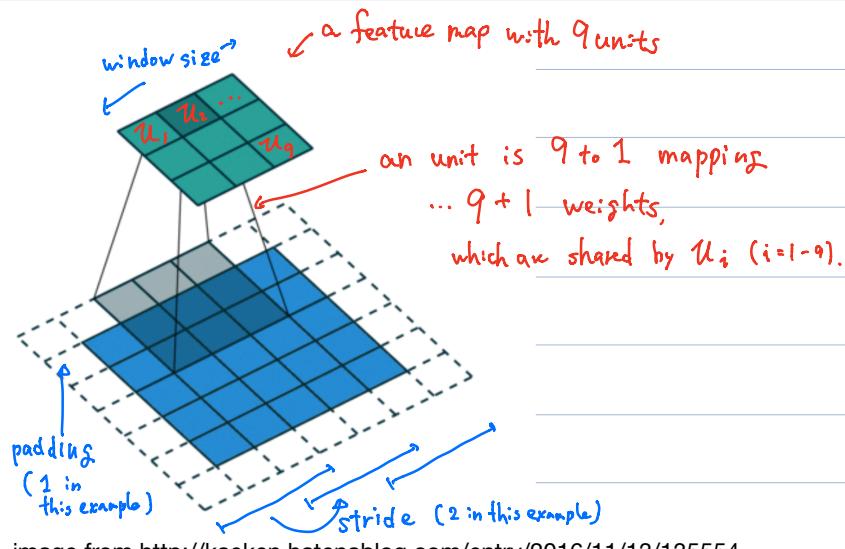


image from <http://kaeken.hatenablog.com/entry/2016/11/13/135554>

a subsampling layer (pooling layer) :

reduce the data size of each feature map (by MAX, AVG, etc.)

## § 5.5. 7 Soft weight sharing

... regularization to motivate similar weight  
( $\equiv$  prior)

$$p(\mathbf{w}) = \prod_{i=1}^W \left( \sum_{j=1}^M \pi_j N(w_i | \mu_j, \sigma_j^2) \right) \Rightarrow \pi_j \text{ is constrained: } \sum \pi_j = 1.$$

$$\Rightarrow \tilde{E}(\mathbf{w}) = E(\mathbf{w}) + \Omega(\mathbf{w}) \quad \text{where}$$

$$\Omega(\mathbf{w}) = - \sum_i \ln \left( \sum_j \pi_j N(w_i | \mu_j, \sigma_j^2) \right)$$

$\rightarrow$  error minimization w.r.t.  $\{w_i, \pi_j, \mu_j, \sigma_j^2\}$

Defining  $\gamma_j(x) = \frac{\pi_j N(x | \mu_j, \sigma_j^2)}{\sum_k \pi_k N(x | \mu_k, \sigma_k^2)}$ ,

$$\frac{\partial \tilde{E}}{\partial w_i} = \frac{\partial E}{\partial w_i} + \sum_{j=1}^M \frac{w_i - \mu_j}{\sigma_j^2} \gamma_j(w_i) \quad \dots w \rightarrow \mu.$$

$$\frac{\partial \tilde{E}}{\partial \mu_j} = - \sum_i \frac{w_i - \mu_j}{\sigma_j^2} \gamma_j(w_i) \quad \dots \mu \rightarrow w.$$

$$\frac{\partial \tilde{E}}{\partial \sigma_j^2} = \sum_i \left( \frac{1}{\sigma_j^2} - \frac{(w_i - \mu_j)^2}{\sigma_j^4} \right) \gamma_j(w_i) \quad [\text{practically } \gamma_j = \log \sigma_j^2 \text{ is used to guarantee } \sigma_j^2 > 0.]$$

I doubt the sentence "It also has the effect of discouraging pathological solutions in which one or more of the  $\sigma_j$  goes to zero" because this is just a variable replacement. Such an effect might be observed in numerical calculation, which is just because of real-number precision.

Noting that  $\sum \pi_j = 1$  with  $0 \leq \pi_j \leq 1$ ,  $\frac{\partial}{\partial \pi}$  should be calculated with

introducing  $\pi_j := \frac{\exp \alpha_j}{\sum \exp \alpha_k}$ . ( $\Leftrightarrow \alpha_j = \log \pi_j + C$ )

$$\begin{aligned} \frac{\partial \tilde{E}}{\partial \alpha_j} &= \frac{\partial \pi_k}{\partial \alpha_j} \frac{\partial \Omega}{\partial \pi_k} \\ &= (\delta_{jk} - \pi_j) \pi_k \cdot \left( - \sum_i \frac{\gamma_k}{\pi_k} \right) \\ &= \sum_i (\pi_j - \gamma_j(w_i)). \quad \dots \pi \rightarrow \gamma. \end{aligned}$$

Direct calculation gives equivalent result:

$$\frac{\partial \Omega}{\partial \pi_j} = \frac{N_j - N_m}{\sum \pi_k N_k} = \frac{\gamma_j}{\pi_j} - \frac{\gamma_m}{\pi_m} \propto \gamma_j \left( 1 - \pi_1 - \dots - \pi_{m-1} \right) - \gamma_m \pi_j$$

## § 5.6 Mixture density network

Non-Gaussian  $P(t|x)$  often appears, e.g., in "inverse problem"

→ use "mixture model" for  $P(t|x)$ . ... multi-modal distribution

$$\sum_i \pi_i(x) \overset{||}{P}_i(t|x)$$

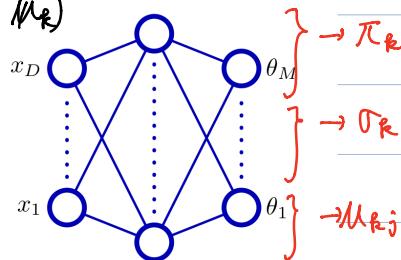
coefficients also depend on  $x$ .

EXAMPLE: Heteroscedastic model

$$P(t|x) = \sum_{k=1}^K \pi_k(x) \mathcal{N}(t | \mu_k(x), \sigma_k^2(x) \mathbb{I})$$

practically we use  $\alpha^\pi$  and  $\alpha^\sigma$  (and  $\alpha^\mu \equiv \mu_k$ )

$$\pi_k := \frac{\exp \alpha_k^\pi}{\sum_l \exp \alpha_l^\pi}$$



$$\sigma_k^2 := \exp \alpha_k^\sigma$$

so that  $\begin{cases} \sum \pi_k = 1, 0 \leq \pi_k \leq 1 \\ \sigma_k^2 > 0 \end{cases}$  are satisfied.

Straightforwardly, with  $\mathcal{N}_{nk} := \mathcal{N}(t_n | \mu_k(x_n), \sigma_k^2(x_n) \mathbb{I})$ ,

$$E(\mathbf{w}) = - \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k(x_n, \mathbf{w}) \mathcal{N}(t_n | \mu_k(x_n, \mathbf{w}), \sigma_k^2(x_n, \mathbf{w}) \mathbb{I}) \right\}$$

$$\frac{\partial E_n}{\partial \alpha_k^\pi} = \pi_k - \gamma_{nk}, \quad \frac{\partial E_n}{\partial \alpha_{nk}^\mu} = \gamma_{nk} \cdot \frac{\mu_{nk} - t_n}{\sigma_{nk}^2}, \quad \frac{\partial E_n}{\partial \alpha_{nk}^\sigma} = \gamma_{nk} \left( L - \frac{\|t_n - \mu_{nk}\|^2}{\sigma_{nk}^2} \right)$$

$$\text{where } L = \dim t, \quad \gamma_{nk} \equiv \gamma_k(t_n | x_n) := \frac{\pi_k \mathcal{N}_{nk}}{\sum_{l=1}^K \pi_l \mathcal{N}_{nl}}.$$

→ We can determine not only

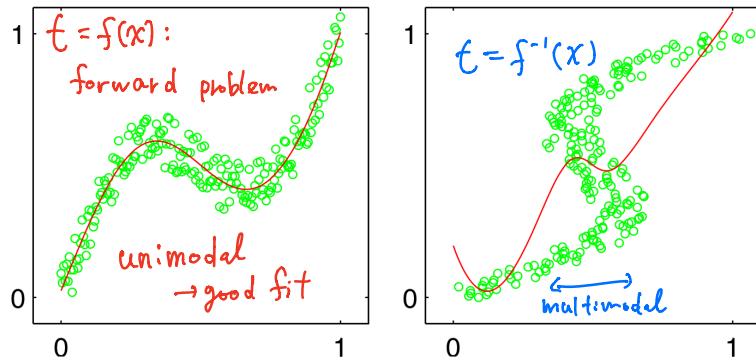
$$\mathbb{E}[t|x] = \int t p(t|x) dt = \sum_{k=1}^K \pi_k(x) \mu_k(x)$$

but also the uncertainty

$$\begin{aligned} s^2(x) &= \mathbb{E} [\|t - \mathbb{E}[t|x]\|^2 | x] \\ &= \sum_{k=1}^K \pi_k(x) \left\{ \sigma_k^2(x) + \left\| \mu_k(x) - \sum_{l=1}^K \pi_l(x) \mu_l(x) \right\|^2 \right\} \end{aligned}$$

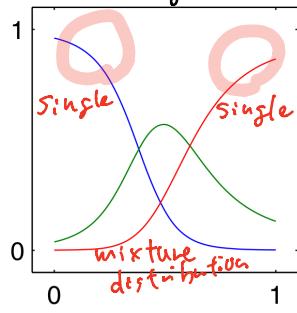
An example: fitting  $f(x) = x + 0.3\sin 2\pi x$

- Simple NN



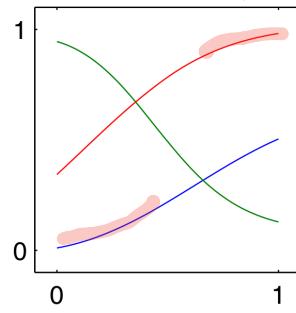
- Mixture density network

the weights  $\pi_k$

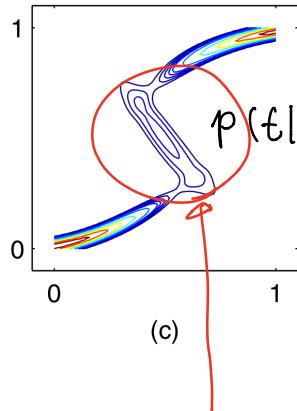


(a)

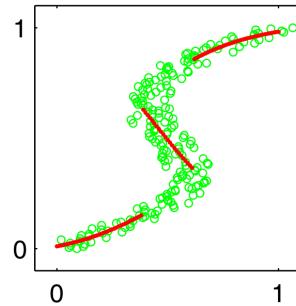
the means  $\mu_k$



"result"  $\mu_k | k \text{ s.t. } \pi_k = \max_i \pi_i$



(c)



(d)

mean  $E[t|x] = \sum \pi_k \mu_k$  is not good for this purpose.

... We should rather use, for example,

the mean  $\mu_k$  of the most probable distribution.

## § 5.7 Bayesian NN.

... no longer analytical due to non-linearity

{ "variational influence" ... § 10.

Laplace approximation

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

$$f'(z_0) = 0 \quad f''(z_0) < 0 \quad A = -\frac{d^2}{dz^2} \ln f(z_0).$$

Given a dataset  $\mathcal{D} = \{(x_1, t_1), \dots, (x_N, t_N)\}$ ,

$$p(w|\mathcal{D}, \alpha, \beta) \propto p(w|\alpha) p(t|x, w, \beta)$$

$$\begin{aligned} & \downarrow \qquad \qquad \downarrow \text{we shall suppose} \qquad \beta^{-1} \equiv \sigma^2 \\ N(w|0, \alpha^{-1} I) \quad & \quad \mathcal{T} N(t_n | \gamma(x_n, w), \beta^{-1}) \\ & \qquad \qquad \qquad \text{NN: non-linear} \end{aligned}$$

Laplace approx. (with  $\alpha$  and  $\beta$  fixed for simplicity)

$$\ln p(w|\mathcal{D}) = -\frac{\alpha}{2} w^T w - \frac{\beta}{2} \sum_n (\gamma(x_n, w) - t_n)^2 + \text{const.}$$

$$\Rightarrow A = -\nabla \nabla \ln p(w|\mathcal{D}) = \alpha I + \beta H; \quad H = \nabla \nabla \frac{1}{2} (y - t)^2.$$

$$p(w|\mathcal{D}) \approx N(w|w_{MAP}, A^{-1}) \quad \text{(back propagation)}$$

Predictive distribution is given by

$$p(t|x, \mathcal{D}) \approx \int dw p(t|x, w, \beta) q(w|\mathcal{D})$$

$$= \int dw \underbrace{N(t | \gamma(x, w), \beta^{-1})}_{\text{again non-linear}} N(w | w_{MAP}, A^{-1})$$

assuming

$$w - w_{MAP} \equiv \Delta w \ll \beta^{-1/2} = \sigma$$

$$\approx N(t | \gamma(x, w_{MAP}) + g^T (w - w_{MAP}), \beta)$$

$$\approx N(t | \gamma(x, w_{MAP}), \sigma^2(x))$$

$$g := \nabla_w \gamma|_{w=w_{MAP}}$$

$$\text{where } \sigma^2(x) = \beta^{-1} + g^T A^{-1} g.$$

## § 5.7.2 Hyperparameters $\alpha$ and $\beta$ can also be optimized :

$$p(t|w, D) = \int d\theta d\alpha d\beta \underbrace{p(t|w, \beta)}_{\text{intrinsic noise}} p(w|D, \alpha, \beta) p(\alpha, \beta|D) \quad \begin{matrix} \text{predictive dist.} \\ \text{parameters from hyperparameter estimation} \end{matrix}$$

$$\ln P(D|\alpha, \beta) = \ln \left[ \int d\theta \frac{\underbrace{p(D|w, \alpha, \beta) p(w|\alpha, \beta)}_{\mathcal{T} \sim N(t_n | y(x_n, w), \beta^{-1})}}{f(w)} \right] \xrightarrow{\text{assume } N(w|\theta, \alpha^{-1}\mathbb{I})} \\ \approx \ln \left[ \int d\theta f(w_{MAP}) \exp \left( -\frac{1}{2} (w - w_{MAP})^T \cancel{A}^{-1} (w - w_{MAP}) \right) \right] \\ f(w) = p(D|w, \alpha, \beta) p(w|\alpha, \beta) \quad \dots \text{numerator of } P(w|D, \alpha, \beta) !$$

$\rightarrow A$  and  $w_{MAP}$  are in the previous page.

$$p(w|D, h) = \frac{p(D|w, h) p(w|h)}{p(D|h)} : \text{Bayes' theorem}$$

$$p(D|h) = \int d\theta p(D|w, h) p(w|h) : \text{marginalization}$$

$$= \ln f(w_{MAP}) - \frac{1}{2} \ln |A| + \frac{W}{2} \log 2\pi \quad \leftarrow \dim w \\ = - \sum_n \frac{\beta}{2} (y(x_n, w_{MAP}) - t_n)^2 + \frac{\alpha}{2} \|w_{MAP}\|^2 \quad + \frac{N}{2} \log \frac{\beta}{2\pi} + \frac{W}{2} \log 2\pi \alpha \\ \equiv E(w_{MAP}) \quad - \frac{1}{2} \ln |A|$$

In evidence framework, we make point estimates for  $\alpha$  and  $\beta$  by maximizing  $\ln P(D|\alpha, \beta)$ .

As § 3.5.1, with eigensystem  $\{\lambda_i, \mathbf{U}_i\}$  of  $\beta H$ ,

$$A = -\nabla \nabla \ln p(w|D) = \alpha I + \beta H ;$$

$$H = \nabla \nabla \frac{1}{2} (y - t)^2 .$$

$$\alpha = \frac{\gamma}{\|w_{MAP}\|^2}, \quad \beta^{-1} = \frac{1}{N-\gamma} \sum_n (y(x_n, w_{MAP}) - t_n)^2, \quad \gamma = \sum_{i=1}^W \frac{\lambda_i}{\alpha + \lambda_i} .$$

and iteratively,  $\{w_{MAP}, \alpha, \beta\}$  is optimized.

Note that multiple (inequivalent) solutions may exist.

I do not understand (and thus doubt) the note that we have ignored the  $\alpha$ -dependence of  $\lambda$ s.  
We can consider  $p$  as a function of  $\{w_{MAP}, \alpha, \beta\}$ , which is valid since we anyway find the minimum iteratively. Then the Hessian has no dependence on  $\alpha$  and  $\beta$ .

$$A = -\nabla \nabla \ln p(w|D) = \alpha I + \beta H ; \quad H = \nabla \nabla \frac{1}{2} (y - t)^2$$

We need  $P(D)$  to compare different models.

$$P(D, \alpha, \beta) = P(D|\alpha, \beta) P(\alpha, \beta) \approx P(D|\hat{\alpha}, \hat{\beta}) \quad (\text{or marginalization:})$$

$$\text{and } P(D) \approx P(D, \alpha, \beta) \times 2^M M ! \quad \text{MacKay 1992c, Bishop 1995a }$$

but note that  $|A|$  is tough to calculate numerically.

§ 9.2.3 For classification.  $\mathcal{X} \mapsto t = \{0, 1\}$

$$P(w|\mathcal{D}, \alpha) \propto P(w|\alpha) P(t|\mathcal{X}, w)$$

$$\rightarrow N(w|0, \alpha^{-1}\mathbb{I}) \prod y_n^{t_n} (1-y_n)^{1-t_n}$$

$$\rightarrow -\ln P(w|\mathcal{D}, \alpha) = E(w) + \text{const.}$$

$$\text{where } E(w) = -\ln P(\mathcal{D}|w) + \frac{\alpha}{2} \|w\|^2.$$

$$\nabla E(w) = \sum_n \frac{y_n - t_n}{y_n(1-y_n)} \nabla y_n + \alpha w \implies w_{\text{MAP}}.$$

$$\nabla \nabla E(w) = \alpha \mathbb{I} + \sum_n \left[ \frac{y_n^2 - 2t_n y_n + t_n}{y_n^2(1-y_n)^2} (\nabla y)(\nabla y) + \frac{y_n - t_n}{y_n(1-y_n)} \nabla \nabla y \right]$$

Hessian H-1

Model discrimination

$$\ln P(\mathcal{D}|\alpha) = \ln \left[ \int d\mathbf{w} P(\mathcal{D}|w) P(w|\alpha) \right]$$

$$\underset{\text{Laplace}}{\approx} \ln f(w_{\text{MAP}}) - \frac{1}{2} \ln |A| + \frac{W}{2} \ln 2\pi \equiv -E(w_{\text{MAP}})$$

$$= \sum_n \left[ t_n \ln y_n + (1-t_n) \ln (1-y_n) \right]_{w_{\text{MAP}}} - \frac{W}{2} \ln \frac{2\pi}{\alpha} - \frac{\alpha}{2} \|w_{\text{MAP}}\|^2$$

$$= -E(w_{\text{MAP}}) - \frac{1}{2} \ln |A| + \frac{W}{2} \ln \alpha - \frac{1}{2} \ln |A| + \frac{W}{2} \ln 2\pi$$

$$\implies \alpha = \frac{\gamma}{\|w_{\text{MAP}}\|^2}, \quad \gamma = \sum_{i=1}^W \frac{\lambda_i}{\alpha + \lambda_i}; \quad \lambda_i \text{ are eigenvalues of } H. \quad (\text{and iterate})$$

Predictive distribution is given by

$$P(t|\mathcal{X}, \mathcal{D}) \approx \int d\mathbf{w} P(t|\mathcal{X}, w) N(w|w_{\text{MAP}}, A^{-1}) \quad \left( \text{or } \approx P(t|\mathcal{X}, w_{\text{MAP}}) \right)$$

Expansion should be done w.r.t.  $a = \sigma^{-1}(y)$ , not  $y$  itself.

$$a(\mathcal{X}, w) \approx a_{\text{MAP}}(\mathcal{X}) + b^T(w - w_{\text{MAP}}) \quad \text{with } b = \nabla a(\mathcal{X}, w_{\text{MAP}}) \text{ and}$$

$$P(t|\mathcal{X}, \mathcal{D}) \approx \int d\mathbf{w} P(t|\mathcal{X}, w) N(w|w_{\text{MAP}}, A^{-1}) : g(a) = \sigma(a)^t (1-\sigma(a))^{1-t}$$

$$= \int d\mathbf{w} da \delta(a - a_{\text{MAP}} - b^T(w - w_{\text{MAP}})) g(a) N(w|w_{\text{MAP}}, A^{-1})$$

$$= \int da g(a) N(a|a_{\text{MAP}}, b^T A^{-1} b)$$

$$\therefore P(t=1|\mathcal{X}, \mathcal{D}) \approx \Gamma(k a_{\text{MAP}}) \quad \text{where } k = \left(1 + \frac{\pi}{8} b^T A^{-1} b\right)^{-1/2}.$$

