# RECAP

$$\text{Data} = \left\{(x_i, t_i)\right\}_{i=1}^{N}$$

(A)

(B)

model
$y = w^{\beta} x$
$y = w^T \phi$

$\downarrow$

Objective func.
(loss)

$\downarrow$

Solve!

model     likelihood     prior
$p(t|x,w)$   $p(w)$
$N(t| v^T\phi, \beta^{-1})$   $N(w| 0, \sigma^{-1} I)$

$\downarrow$

Inference.

MLE,  MAP   (Bayesian)

$\downarrow$

Solve!

Decision!
for fun & profit!



$\lambda \uparrow$   $\lambda \downarrow$

error

test error

Training Error

A       model Complexity      B

$\boxed{\frac{\lambda}{2}} w^T w$

A   underfitting

B   overfitting

$w^{(1)} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$

$[w_1 \; w_2] \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$
$= w_1^2 + w_2^2 \leq c$
$w_1^2 + w_2^2 = c$

$$\min_{w} \frac{1}{2}(X_w - t)^2 + \underbrace{\frac{\lambda}{2} \|w\|_2^2}_{w^T w}$$

$\Downarrow$

$$\min_{w} \frac{1}{2}(X_w - t)^2$$
$$\text{s.t} \quad w^T w \leq c$$

Practice.

|  | 80% | 20% |  |
|---|---|---|---|
| D = | Training | Validation | test |

$\lambda = \{10^{-5}, 10^{-4} \cdots\}$

Validation

X-Validation

| D = | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|

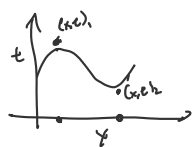Bias - Variance Decomposition / Trade-Off.

Key idea!

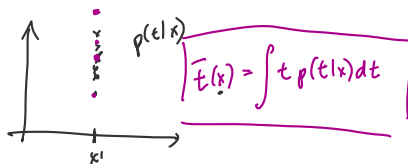$$\boxed{\text{Generalization Error} = \text{Bias} + \text{Variance} + \text{Noise}.}$$
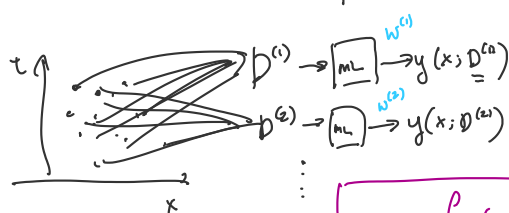
**SETUP**

Distributions

① $p(x,t) = p(x)p(t|x)$



samples drawn i.i.d.
independent or identically distributed.

$$\bar{t}(x) = \int t\, p(t|x)\, dt$$

② Dataset Distribution $p^N$ · $p(D)$



$D^{(1)} \to \boxed{ML} \xrightarrow{w^{(1)}} y(x; D^{(1)})$

$D^{(2)} \to \boxed{ML} \xrightarrow{w^{(2)}} y(x; D^{(2)})$

$$p(D) = \prod_{n=1}^{N} p(x_n, t_n)$$

samples, with replacement

$[x_1, x_2, x_3, x_4 \ldots, x_{10}]$

$D^{(1)} = \{x_2, x_7\}$

$D^{(2)} = \{x_1, x_7\}$

$p(D^{(1)}) = p(x_2) p(x_7)$
$\qquad = \frac{1}{10} \cdot \frac{1}{10} = \frac{1}{100}$

$$\bar{y}(x) = \int_D y(x; D)\, p(D)\, dD$$

$$\int_w y(x; w)\, p(w)\, dw$$

$$= \mathbb{E}_D[y(x; D)]$$

Gen. Error

$$\mathbb{E}_{x, t, D}\left[ \mathcal{L}(x, t, D) \right]$$

$$= \mathbb{E}_{x, t, D}\left[ (y(x; D) - t)^2 \right]$$

$$= \mathbb{E}_{x, t, D}\left[ \underbrace{(y(x; D) - \bar{y}(x)}_{a} + \underbrace{\bar{y}(x) - t)}_{b})^2 \right] \qquad a^2 + b^2 + 2ab$$

$$= \boxed{\mathbb{E}\left[ \underline{(y(x; D) - \bar{y}(x))^2} \right]} + \mathbb{E}\left[ (\bar{y}(x) - t)^2 \right] + \mathbb{E}\left[ 2(y(x; D) - \bar{y}(x))(\bar{y}(x) - t) \right]$$

Claim: $= 0$

$$\mathbb{E}_{x,t}\left[ \mathbb{E}_D\left[ (y(x; D) - \bar{y}(x))(\bar{y}(x) - t) \right] \right]$$

$$\mathbb{E}_{x,t}\left[ (\underbrace{\mathbb{E}_D[y(x; D)]}_{\bar{y}(x)} - \bar{y}(x))(\bar{y}(x) - t) \right]$$

$\underbrace{\qquad\qquad}_{=0} \qquad\qquad =0$

$$\mathbb{E}_{x, t, D}\left[ (\bar{y}(x) - t)^2 \right]$$

$$\mathbb{E}_{x, t, D}\left[ (\bar{y}(x) - \bar{t}(x)) + (\bar{t}(x) - t))^2 \right]$$

$$\boxed{\mathbb{E}_{x, t, D}\left[ (\bar{y}(x) - \bar{t}(x))^2 \right]} + \mathbb{E}_{x, t, D}\left[ (\bar{t}(x) - t)^2 \right] + \mathbb{E}\left[ 2(\bar{y}(x) - \bar{t}(x))(\bar{t}(x) - t) \right]$$

Claim: $= 0$

Exercise $\qquad p(t|x)$

$$\mathbb{E}_{x, t, D}\left[ (y(x) - t)^2 \right] = \mathbb{E}_{x, t, D}\left[ (\bar{y}(x) - \bar{t}(x))^2 \right] + \mathbb{E}_{x, t, D}\left[ (y(x; D) - \bar{y}(x))^2 \right] + \mathbb{E}_{x, t, D}\left[ (\bar{t}(x) - t)^2 \right]$$
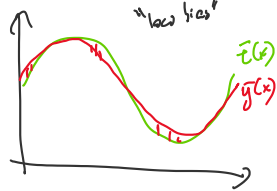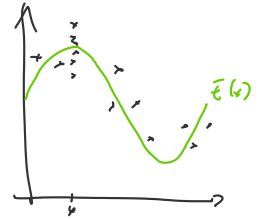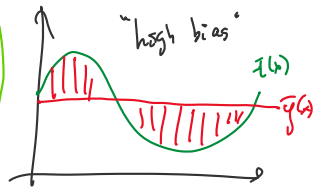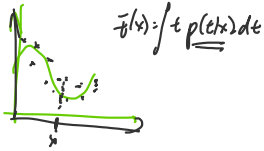
"Variance" $\qquad\qquad$ "noise"

**example:**

$h(x) = \sin(2\pi x)$

$t(x) = h(x) + \varepsilon$

$\sim N(h(x), \sigma_n^2) \quad \sim N(0, \sigma_n^2)$

$\bar{t}(x) = h(x)$

"high bias"   $\bar{t}(x)$   $g(x)$

"low variance"   $y(x; D^1)$   $y(x; D^2)$   $g(x)$   $y(x; D^3)$

$\bar{t}(x)$

$\bar{t}(x) = \int t \, p(t|x) \, dt$

"low bias"   $\bar{t}(x)$   $g(x)$

"high variance"   $y(x; D^1)$   $\bar{y}(x)$   $y(x; D^2)$

**Putting it together**

Regime A
high bias

Regime B
high variance

Test error

variance

error

Bias

Training error

model complexity

error

regime B
high variance

regime A
high bias

Test/Gen.

Train

error Threshold.

model
$y = w^T x$

# Training Pts.

**Bayesian Approach**

Recall: $p(t|x, w) = N(t | w^T \phi, \beta^{-1})$

say, we have $p(w|D)$   "posterior weight distribution"

Want $p(t|x)$

(sum rule)

$\sum_w p(t|x,w) p(w|D)$

$p(t|x) = \int_w p(t|x,w) \, p(w|D) \, dw$

$= \int_w N(t | w^T \phi, \beta^{-1}) \, \boxed{N(w | m_N, S_N)} \, dw \quad \leftarrow ?$

$= N\left(t \,\middle|\, m_N^T \phi(x), \; \sigma_N^2(x)\right)$

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

learning $p(w|D)$

likelihood   prior

$p(D|w) \; p(w)$

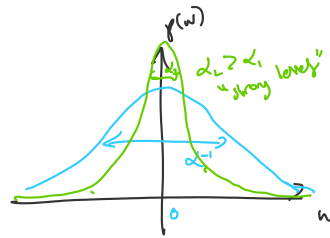$$p(w|D) = \frac{\boxed{p(D|w)p(w)}}{p(D)} \quad \text{(Bayes rule)}$$

$$p(D) \xleftarrow{} \text{evidence.} = \int_w p(D|w)p(w)\,dw$$

<u>Setup</u>

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

$$\begin{bmatrix} \frac{1}{\alpha} & & 0 \\ & \frac{1}{\alpha} & \\ & & \ddots \\ 0 & & \frac{1}{\alpha} \end{bmatrix}$$


$p(w)$, $\alpha_2 > \alpha_1$ "strong lovers", $\alpha^{-1}$, $0$, $w$

likelihood: $\quad p(D|w) = \prod_{n=1}^{N} N(t_n\,|\,\boxed{w^T \phi_n}, \beta^{-1}) \xleftarrow{} \text{assume known.}$

$$= N(t\,|\,\Phi w,\; \beta^{-1} I)$$

$t_n = w^T \phi_n + \varepsilon$

$\varepsilon \sim N(0, \beta^{-1})$

<u>Exercise.</u>

$$\begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{bmatrix} \quad \begin{bmatrix} -\phi_{(x_1)}- \\ -\phi_{(x_2)}- \\ \vdots \\ -\phi_{(x_N)}- \end{bmatrix} \quad \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_{m-1} \end{bmatrix}$$

---

Revision $\quad$ (See 2.3 of PRML) $\quad$ Multivariate Gaussians.


$p(x)$, $\sigma^2$, $x$

$p(x)$ $\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$, $\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$, $x_2$, $x_1$

Covariance matrix $\Sigma = \begin{bmatrix} \sigma_{11}^2 & \sigma_{12}^2 \\ \sigma_{21}^2 & \sigma_{22}^2 \end{bmatrix}$

$x_2$, $x_1$

$$x \sim N(x\,|\,\mu, \Sigma)$$

$$N(x\,|\,\mu, \Sigma) = \underbrace{\frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}}}_{\text{normalizer.}} \exp\left[ -\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) \right]$$

$$-\frac{(x-\mu)^2}{2\sigma_x^2}$$

Covariance Matrices

$$\begin{bmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{bmatrix} \qquad \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix} \qquad \begin{bmatrix} \sigma_{11}^2 & \sigma_{12}^2 \\ \sigma_{21}^2 & \sigma_{22}^2 \end{bmatrix}$$

"spherical" $\qquad$ "diagonal" $\qquad$ "full"


$x_2$, $x_1$ $\qquad$ $x_2$, $x_1$ $\qquad$ $x_2$, $x_1$

---

(2.116) in PRML

<u>Marginal</u> & <u>Conditional Gaussian Distributions</u>

Given a marginal Gaussian for $w$
and a conditional Gaussian for $t$ given $w$.

$$p(w) = N(w\,|\,\mu, \Lambda^{-1}) \quad \text{"marginal"}$$

$$p(t|w) = N(t\,|\,Aw+b,\; L^{-1}) \quad \text{"conditional"}$$

Then

the marginal
$$p(t) = N(t\,|\,A\mu + b,\; L^{-1} + A\Lambda^{-1} A^T)$$

$$p(w|t) = N\left(w \mid \Sigma\left(A^T L (t-b) + \Lambda \mu\right), \Sigma\right)$$

$$\text{where} \quad \Sigma = \left(\Lambda + A^T L A\right)^{-1}$$

$$p(w|D) \Rightarrow \frac{N(t \mid \Phi w, \beta^{-1} I)\, N(w \mid 0, \alpha^{-1} I)}{p(v) = \int_w N(t \mid \Phi w, \beta^{-1} I)\, N(w \mid 0, \alpha^{-1} I)\, dw}$$

Posterior:

$$p(w|t) = N(w \mid M_N, S_N)$$

$$\Rightarrow M_N = \beta S_N \Phi^T t$$

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

$$M_N = \beta \left(\beta \Phi^T \Phi + \alpha I\right)^{-1} \Phi^T t$$

$$= \underbrace{\left(\Phi^T \Phi + \frac{\alpha}{\beta} I\right)^{-1} \Phi^T}_{\text{pseudo inverse}} t$$

(MAP)