

# General decomposition of the generalization error into Variance - Bias - Noise

Consider :

- 1.) Sample space  $X$
- 2.) Target space  $y$
- 3.) training data  $D = \{(X_1, Y_1), \dots, (X_N, Y_N)\}$   
as a random variable  
on  $(X \times Y)^N$
- 4.) test data point  $(X, Y) \in X \times Y$   
 $D$  and  $(X, Y)$  sampled from  
"true" distribution  $P_0$

5.) loss function :

$$L : \mathcal{Y} \times \mathcal{P}(\mathcal{Y}) \rightarrow \mathbb{R} \cup \{\pm\infty\}$$

$\nearrow$   
space of all probability  
distributions / densities on  $\mathcal{Y}$

- s.t. a)  $L$  is convex in 2nd argument.  
b)  $\{y \sim p(y) \Rightarrow p \in \arg \min_q E_p [L(y, q)]\}$

6.) Any statistical model class  
 $\mathcal{Q} = \{q(y|x)\}$

and learning algorithm  $\delta$   
(e.g. MLE, MAP, full Bayesian)

together determine via learning  
on  $D$  a map ;  $q_D$  :

$$q(y|x, D) : (\mathcal{X} \times \mathcal{Y})^N \times \mathcal{X} \rightarrow \mathcal{P}(\mathcal{Y})$$

training test x prediction  
data data input for  $y$

$$7.) \text{ Put : } \bar{q}(Y|X) := E_q[q(Y|X, D)] \\ = \int q(Y|X, D) p_0(D|x) dD$$

the mean prediction probability  
averaged over all data sets.

(i.e. sampling several  $D_1, D_2, \dots \sim p_0$   
running on each the learning algorithm  
A and afterwards averaging)

$$8.) \text{ Let } q_0(Y|X) \in \underset{q(Y|X)}{\operatorname{arg\min}} E_q[L(Y, q(Y|X))|X]$$

the optimal prediction probability  
which usually is  $p_0(Y|X)$ .  $\leftarrow$  "true"  
(depending on L)

We then define:

### 9.) Generalization error:

$$\mathcal{E}_L(q_{\mathcal{D}}|X) := \mathbb{E}_0 [L(Y, q(Y|X, \mathcal{D})) | X]$$

### 10.) Variance:

Jennsen

$$\text{Var}_L(q_{\mathcal{D}}|X) := \mathbb{E}_0 [L(Y, q(Y|X, \mathcal{D})) | X] - \mathbb{E}_0 [L(Y, \bar{q}(Y|X)) | X] \geq 0$$

### 11.) Bias:

$p_0 \in \arg\min$

$$\text{Bias}_L(q_{\mathcal{D}}|X) := \mathbb{E}_0 [L(Y, \bar{q}(Y|X)) | X] - \mathbb{E}_0 [L(Y, p_0(Y|X)) | X] \geq 0$$

### 12.) Noise:

$$\text{Noise}_L(q_{\mathcal{D}}|X) := \mathbb{E}_0 [L(Y, p_0(Y|X)) | X]$$

↳ independent of  $\mathcal{Q}$  and  $\mathcal{A}$ .  
↳ non-reducible noise.

Then we have the decomposition :

$$\mathcal{E}_L = \text{Var}_L + \text{Bias}_L + \text{Noise}_L$$

$$\mathcal{E}_L \geq \text{Noise}_L$$

Example :  $L(y, g) := (y - E_{y|g}[y])^2$   
"square-loss"

$$\text{Var}_L = E_0 \left[ \left( E_{g(Y|X,D)}(Y|X) - E_{\bar{g}(Y|X)}(Y|X) \right)^2 \middle| X \right]$$

$$\text{Bias}_L = \left( E_{\bar{g}}(Y|X) - E_0(Y|X) \right)^2$$

$$\text{Noise}_L = \text{Var}_0(Y|X)$$

$$\mathcal{E}_L = E_0 \left[ (Y - E_{g(Y|X,D)}(Y|X))^2 \mid X \right]$$

⇒ Recover old formula

Check:

$$\text{Var}_L = \mathbb{E}_0 [ L(Y, q(Y|X, D)) | X ]$$

$$- \mathbb{E}_0 [ L(Y, \bar{q}(Y|X)) | X ]$$

$$= \mathbb{E}_0 [ (Y - \mathbb{E}_q[Y|X, D])^2 | X ]$$

$$- \mathbb{E}_0 [ (Y - \mathbb{E}_{\bar{q}}[Y|X])^2 | X ]$$

$$= \mathbb{E}_0 [ Y^2 - 2 \cdot Y \cdot \mathbb{E}_q[Y|X, D] + \mathbb{E}_q[Y|X, D]^2$$

$$(-Y^2 + 2 \cdot Y \cdot \mathbb{E}_{\bar{q}}[Y|X] - \mathbb{E}_{\bar{q}}[Y|X]^2) | X ]$$

$$= \mathbb{E}_0 [ \mathbb{E}_q[Y|X, D]^2 - 2 \cdot Y \cdot \mathbb{E}_q[Y|X, D] | X ]$$

$$+ 2 \mathbb{E}_0[Y|X] \cdot \mathbb{E}_{\bar{q}}[Y|X] - \mathbb{E}_{\bar{q}}[Y|X]^2 \uparrow p_0(Y, D|X)$$

$Y \perp\!\!\!\perp D|X$

$$= \mathbb{E}_0 [ \mathbb{E}_q[Y|X, D]^2 - 2 \mathbb{E}_0[Y|X] \cdot \mathbb{E}_{\bar{q}}[Y|X] | X ] = p_0(Y|X) \cdot p_0(D|X)$$

$$+ 2 \mathbb{E}_0[Y|X] \cdot \mathbb{E}_{\bar{q}}[Y|X] - \mathbb{E}_{\bar{q}}[Y|X]^2$$

$$= \mathbb{E}_0 [ (\mathbb{E}_q[Y|X, D] - \mathbb{E}_{\bar{q}}[Y|X])^2 | X ]$$

$$\text{Bias}_L = \mathbb{E}_0[(Y - \mathbb{E}_{\bar{q}}(Y|X))^2 | X]$$

$$- \mathbb{E}_0[(Y - \mathbb{E}_0(Y|X))^2 | X]$$

$$= \mathbb{E}_0[\cancel{Y^2} - 2Y \cdot \mathbb{E}_{\bar{q}}(Y|X) + \mathbb{E}_{\bar{q}}(Y|X)^2$$

$$- \cancel{Y^2} + 2Y \cdot \mathbb{E}_0(Y|X) + \mathbb{E}_0(Y|X)^2 | X]$$

$$= \mathbb{E}_{\bar{q}}(Y|X)^2 - 2\mathbb{E}_0(Y|X) \cdot \mathbb{E}_{\bar{q}}(Y|X) + \mathbb{E}_0(Y|X)^2$$

$$= (\mathbb{E}_{\bar{q}}(Y|X) - \mathbb{E}_0(Y|X))^2$$

$$\text{Noise}_L = \mathbb{E}_0[(Y - \mathbb{E}_0(Y|X))^2 | X]$$

$$= \text{Var}_0(Y|X)$$

Example :  $L(y, q) := -\log q(y)$

a log-loss

$$\text{Var}_L = \mathbb{E}_0 \left[ \log \frac{\bar{q}(y|x)}{q(y|x,\theta)} \right] \geq 0$$

$$\text{Bias}_L = KL(p_0(y|x) \| \bar{q}(y|x)) \geq 0$$

$$\text{Noise}_L = H(p_0(y|x))$$

$$\mathcal{E}_L = CE[p_0(y|x) \| q(y|x,\theta)]$$

or generalization error beyond noise:

$$KL(p_0(y|x) \| q(y|x,\theta))$$

$$= \mathbb{E}_{y \sim p_0(y|x)} \left[ \log \frac{\bar{q}(y|x)}{q(y|x,\theta)} \right] + \underbrace{KL(p_0(y|x) \| \bar{q}(y|x))}_{\text{Var}_L \geq 0} + \underbrace{\text{Bias}_L}_{\geq 0}$$