

Loss: Squared error loss

Input:  $x$

Prediction:  $y$

Targets: conditional distribution  $P_D(t|x)$

Minimizing the expected loss:

$$\begin{aligned} \mathbb{E}_{P_D} \left[ \underbrace{(y-t)^2}_{\text{squared error loss}} \mid x \right] &= \mathbb{E} \left[ y^2 - 2yt + t^2 \mid x \right] \\ &= y^2 - 2y \mathbb{E}[t|x] + \mathbb{E}[t^2|x] \\ &= y^2 - 2y \mathbb{E}[t|x] + \mathbb{E}[t|x]^2 + \text{Var}[t|x] \\ &= \underbrace{(y - \mathbb{E}[t|x])^2}_{y^*} + \underbrace{\text{Var}[t|x]}_{\text{Bayes error}} \end{aligned}$$

$\text{Var}[z] = \mathbb{E}[z^2] - \mathbb{E}[z]^2$

$\left\{ \begin{array}{l} (y - y^*)^2 \geq 0 \\ \text{Var}[t|x] \text{ does not depend on } y \end{array} \right.$

Bayes error: best possible generalization error we can achieve (if we model data perfectly)  $\approx$  best risk we can achieve

Treat  $y$  as a random variable

- Experiment:
1. Sample a training set from  $P_D$
  2. Train a model (e.g., neural network)
  3. Compute predictions on  $x$

(for simplicity, omit the dependence on  $x$ )

$$\begin{aligned} \mathbb{E}[(y-t)^2] &= \mathbb{E}[(y-y^*)^2] + \text{Var}(t) \\ &= \mathbb{E}[y^2 - 2yy^* + y^{*2}] + \text{Var}(t) \\ &= \mathbb{E}[y^2] - 2y^* \mathbb{E}[y] + y^{*2} + \text{Var}(t) \\ &= y^{*2} + \underbrace{\mathbb{E}[y^2] + \text{Var}(y)}_{\text{variance}} - 2y^* \mathbb{E}[y] + \text{Var}(t) \\ &= y^{*2} - 2y^* \mathbb{E}[y] + \mathbb{E}[y]^2 + \text{Var}(y) + \text{Var}(t) \end{aligned}$$

**Bias-variance decomposition**

$$\mathbb{E}[(y-t)^2] = \underbrace{(y^* - \mathbb{E}[y])^2}_{\text{bias}} + \underbrace{\text{Var}(y)}_{\text{variance}} + \underbrace{\text{Var}(t)}_{\text{Bayes error}}$$

bias: how far our model's average prediction  
 variance: of predictions explained by the choice of training set (how much does the model overfit to the train set)

Visualization in the output space (2 test examples)

