

Week 4: Model selection & validation

Hyperparameters

Ridge Regression with poly. expansion

$$\min_w \frac{1}{N} \sum_{n=1}^N (y_n - x_n^T w)^2 + \lambda \|w\|_2^2 \rightarrow \lambda, \alpha$$

- Generally, the data $\sim D$ with unknown

True err.

$$L_D(f) = E_D[\ell(y, f(x))]$$

Empirical err.

$$L_S(f) = \frac{1}{|S|} \sum_{(x_n, y_n) \in S} \ell(y_n, f(x_n))$$

$$f_S = A(S)$$

Learn Alg.

$$L_S(f_S) = \frac{1}{|S|} \sum_{(x_n, y_n) \in S} \ell(y_n, f_S(x_n))$$

Train-Test split

$$S = S_{\text{Train}} \cup S_{\text{Test}}, f_S = A(S_{\text{Train}})$$

$$L_{S_{\text{Test}}}(f_S) = \frac{1}{|S_{\text{Test}}|} \sum_{(x_n, y_n) \in S_{\text{Test}}} \ell(y_n, f_S(x_n))$$

The true err is:

$$L_D(f) = E_{(y, x) \sim D} [\ell(y, f(x))]$$

empirical err is unbiased:

$$|L_D(f) - L_{S_{\text{Test}}}(f)|$$

generalization err:

$$L_D(f) = E_{x_{\text{test}} \sim D} [L_{S_{\text{Test}}}(f)]$$

Variation:

The loss decreases in

$$\mathcal{O}(1/\sqrt{s_{\text{true}}})$$

\rightarrow More data \rightarrow More confidence

$$\mathbb{P}[|L_D(f) - L_{S_{\text{true}}}(f)| \geq \sqrt{\frac{(b-a)^2 \ln(2/\delta)}{2s_{\text{true}}}}] \leq \delta$$

Given $L(\cdot, \cdot)$ is bounded
in $[a, b]$

Chebyshev Bound:

$\theta_1, \dots, \theta_N$, iid with mean $\mathbb{E}[\theta]$ and range $[a, b]$
Then for any $\epsilon > 0$,

$$\mathbb{P}\left[\left|\frac{1}{N} \sum_{n=1}^N \theta_n - \mathbb{E}[\theta]\right| \geq \epsilon\right] \leq 2e^{-\frac{N\epsilon^2}{(b-a)^2}}$$

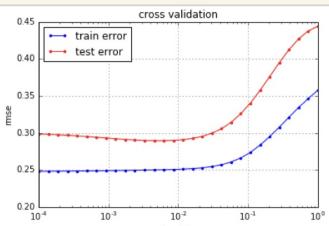
• Equating $2e^{-\frac{1}{s_{\text{true}}} \frac{\epsilon^2}{(b-a)^2}} = \delta$

$$\epsilon = \sqrt{\frac{(b-a)^2 \ln(2/\delta)}{2s_{\text{true}}}}$$

$$\mathbb{P}[L_D(f) > L_S(f) + \sqrt{\frac{(b-a)^2 \ln(2/\delta)}{2s_{\text{true}}}}] \leq \delta$$

• Probabilistic upper bound on the true risk since both $L_S(f)$ as well as the error term can be computed with the doctor.

Model Selection:



Select hyperparameters of our model.

(e.g. in ridge regression)

$S = S_{\text{Train}} \cup S_{\text{Test}}$ and we choose $S_{\text{Test}} \cap S_{\text{Train}} \sim D$

$\lambda_k \mid k = 1 \dots K$

- run the learning algorithm K times on S_{Train} .
- $\rightarrow f_{S_{\text{Train}}, k}$ and compute Training err and Test err.

\rightarrow we select the best model $(f_{S_{\text{Test}} \cap S_{\text{Train}}, k})$

- Similarly as before we have:

The maximum deviation for all K candidates is bounded

$$P\left[\max_k |L_D(f_k) - L_{S_{\text{Test}}}(f_k)| > \sqrt{\frac{(b-a)^2 \ln(2K/\delta)}{2|S_{\text{Test}}|}}\right] \leq \delta$$

- err decreases as $O(\sqrt{|S_{\text{Test}}|})$

and so we have K hyperparameters, can only go up by a small factor $\sim \sqrt{\ln(K)}$

$$\begin{aligned} K^* &= \arg \min_k L_D(f_k) \\ K &= \arg \min_k L_{S_{\text{Test}}}(f_k) \end{aligned}$$

$$P\left[L_D(f_K) > L_D(f_{K^*}) + \sqrt{\frac{(b-a)^2 \ln(2K/\delta)}{2|S_{\text{Test}}|}}\right] \leq \delta$$

In words, if we choose the "best" function according to the empirical risk then its true risk is not too far away from the true risk of the optimal choice.

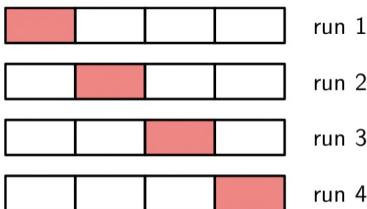
Cross Validation :

Splitting the data into 2 subsets, Train and Test is not the most efficient way to use the data.

K-fold cross-validation is a popular variant.

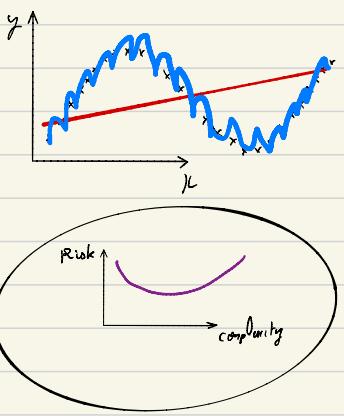
- Randomly partitions into K groups.
- Train K times, each time leaving exactly one of the K groups for testing.
- Average the K results.

Note: Have used all data for training, and all data for testing, and used each data point the same number of times.

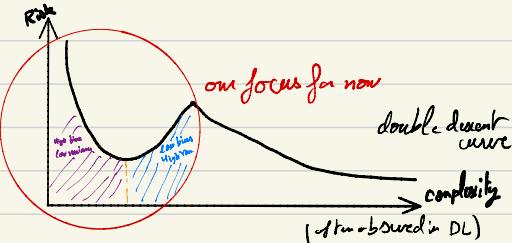


Cross-validation returns an unbiased estimate of the *generalization error* and its variance.

Bias - Variance decomposition



Important note



Data model:

$$y = f(x) + \epsilon$$

true model noise $\epsilon \sim D_\epsilon \text{ iid}, \perp\!\!\!\perp D_x$

$$\mathbb{E}[\epsilon] = 0$$

In general, the model will not be realizable i.e. f is not in our model class (cannot recover exactly)

$$\mathbb{E}_{D_x}[y - f_{S_{T_m}}(x)]^2 \quad (f_{S_{T_m}} = A(S_{T_m}))$$

$$\begin{aligned} & \text{Look at one single point } x_0: (y - f_{S_{T_m}}(x_0))^2 = (f(x_0) + \epsilon - f_{S_{T_m}}(x_0))^2 \\ & \mathbb{E}_{\substack{\epsilon \sim D_\epsilon \\ S_{T_m}}} [(y - f_{S_{T_m}}(x_0))^2] = \mathbb{E}[(f(x_0) + \epsilon - f_{S_{T_m}}(x_0))^2] \\ & = \mathbb{E}_f[(f(x_0) - f_{S_{T_m}}(x_0))^2] + \mathbb{E}[\epsilon^2] + \text{Cross term} \quad \left| \begin{array}{l} \text{Cross term: } \mathbb{E}[2(f(x_0) - f_{S_{T_m}}(x_0))\epsilon] \\ \mathbb{E}[D_\epsilon] = 0 \end{array} \right. \\ & = \mathbb{V}\text{ar}(\epsilon) (\mathbb{E}[\epsilon]^2 = 0) \end{aligned}$$

$$\Rightarrow \mathbb{E}[(y - f_{S_{T_m}}(x_0))^2] = \mathbb{E}[(f(x_0) - f_{S_{T_m}}(x_0))^2] + \mathbb{V}\text{ar}(\epsilon)$$

$$\begin{aligned} & \mathbb{E}[(f(x_0) - f_{S_{T_m}}(x_0))^2] = \mathbb{E}_f[(f(x_0) - \mathbb{E}_{S_{T_m}}[f_{S_{T_m}}(x_0)] + \mathbb{E}_{S_{T_m}}[f_{S_{T_m}}(x_0)] - f_{S_{T_m}}(x_0))^2] \\ & = (f(x_0) - \mathbb{E}_{S_{T_m}}[f_{S_{T_m}}(x_0)])^2 + \mathbb{E}_{S_{T_m}}[(f_{S_{T_m}}(x_0) - \mathbb{E}_{S_{T_m}}[f_{S_{T_m}}(x_0)])^2] + \text{C.T.} \end{aligned}$$

$$\text{C.T. } \mathbb{E}_{S_{T_m}}[(f(x_0) - \mathbb{E}_{S_{T_m}}[f_{S_{T_m}}(x_0)])^2] = \mathbb{E}_{S_{T_m}}[(f(x_0) - f_{S_{T_m}}(x_0))^2] = 0$$

$$\begin{aligned} & \mathbb{E}_{\substack{\epsilon \sim D_\epsilon \\ S_{T_m}}} [(y - f_{S_{T_m}}(x_0))^2] = \mathbb{V}\text{ar}(\epsilon) + (f(x_0) - \mathbb{E}_{S_{T_m}}[f_{S_{T_m}}(x_0)])^2 + \mathbb{E}_{S_{T_m}}[(f_{S_{T_m}}(x_0) - \mathbb{E}_{S_{T_m}}[f_{S_{T_m}}(x_0)])^2] \end{aligned}$$

Bias **Variance**

\Rightarrow Decomposition in 3 parts: Noise term, Bias term, Variance term (over data)

Bias term: $f(x_0)$: actual fct. $\mathbb{E}_{S_{T_m}}[f_{S_{T_m}}(x_0)]$: expectation of the model. (How precise)

\rightarrow If the complexity is high \rightarrow low bias | if complexity is low \rightarrow High bias

Variance term: $f_{S_{T_m}}(x_0)$: mod. $\mathbb{E}_{S_{T_m}}[f_{S_{T_m}}(x_0)]$: average mod. (Exploit) (How consistent)

\rightarrow If the complexity \nearrow Then Variance \nearrow | if complexity \downarrow , Var. \downarrow

Training Risk:



True Risk: Mean-Neg (5000): (Pmt. case)

