

— Any Modification made to learning algorithm to reduce Test error, but not necessarily training error, is termed as regularization.)

— A particular approach to regularization is Weight decay.

in fact. $E^{(train)}$ can increase.

$$\left\{ \min_{\underline{W}} J(\underline{W}) = \sum_{n=1}^N [y^{(n)} - t^{(n)}]^2 \right. \quad \left. \text{subject to } \|\underline{W}\|_2 \leq \rho \right\} \quad L_2\text{-norm}$$

equivalently

$$\min_{\underline{W}} J(\underline{W}) = \sum_{n=1}^N [y^{(n)} - t^{(n)}]^2 + \lambda \|\underline{W}\|_2^2 \quad \text{for } \lambda \geq 0$$

L_2 -norm
→ Ridge Regression

L_1 -norm
→ Lasso Regression
it may lead to more sparse solns.

$$\nabla_{\underline{W}} J(\underline{W}) = 2\underline{W}^T \underline{X}^T \underline{X} - 2\underline{t}^T \underline{X} + 2\lambda \underline{W}^T = 0$$

$$(\underline{X}^T \underline{X} + \lambda \underline{I}) \underline{W} = \underline{X}^T \underline{t}$$

Where \underline{I} is $(D+1) \times (D+1)$ identity matrix

$$\underline{W} = (\underline{X}^T \underline{X} + \lambda \underline{I})^{-1} \underline{X}^T \underline{t}$$

Regularization with nonlinear models

$$\min_{\underline{W}} J(\underline{W}) = \sum_{n=1}^N [\underline{W}^T \underline{\Phi}(x^{(n)}) - t^{(n)}]^2 + \lambda \underline{W}^T \underline{W}$$

$$\underline{W} = (\underline{\Phi}^T \underline{\Phi} + \lambda \underline{I})^{-1} \underline{\Phi}^T \underline{t}$$

$(M+1) \times 1$

$(M+1) \times (M+1)$

$(M+1) \times (M+1)$

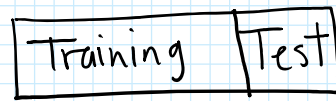
$N \times 1$

$\phi: N \times (M+1)$

$$(M+1) \times 1 + \overbrace{(M+1) \times (M+1)} + (M+1) \times (M+1)$$

① Hold-out Method

→ randomize data set



→ 80% 20%

example or training set

– \underline{W} is generated using training set & Test set is used with \underline{W} to compute $E^{(test)}$

Problems

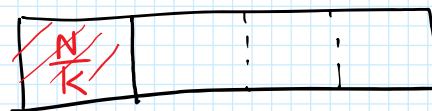
→ for small N (sparse datasets) this is not an efficient way to use data.

→ Unlucky & produce unfortunate splits

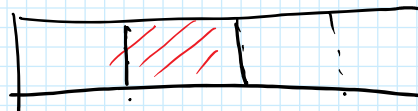
② K-fold Cross-Validation

– K experiments, in each we use $K-1$ folds for training & remaining fold for testing

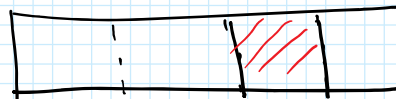
$K=4$



used for testing



$$E_{\bar{K}}^{(test)} = \frac{1}{K} \sum_{i=1}^K E_i^{(test)}$$



– All examples are eventually used for both Testing & training.

③ Leave-one-out is K -fold C.V. with $K=N$

Choice of k depends on N

- for sparse datasets, k is to be large

- for large N , $k=3$ is sufficient

↳ in general $3 \leq k \leq 10$

④ Bootstrap: resampling technique w/ replacement

example:
 $N=5$

$x^{(1)} \quad x^{(2)} \quad x^{(3)} \quad x^{(4)} \quad x^{(5)}$

Exp. 1

$x^{(1)} \quad x^{(3)} \quad x^{(3)} \quad x^{(3)} \quad x^{(5)}$ $x^{(2)} \quad x^{(4)}$ ←
Train. Set test set

Exp. 2

$x^{(1)} \quad x^{(2)} \quad x^{(3)} \quad x^{(4)} \quad x^{(5)}$ $x^{(4)}$

⋮

Exp. K

$$E^{(test)} = \frac{1}{K} \sum_{i=1}^K E_i^{(test)}$$

Hyperparameter Selection

Hyperparameter: Any setting (parameter) that is not adjusted during learning.

e.g., β , λ , M

1- Data is divided into 3 disjoint sets

Training Set	Validation Test	Test Set
60%	20%	20%

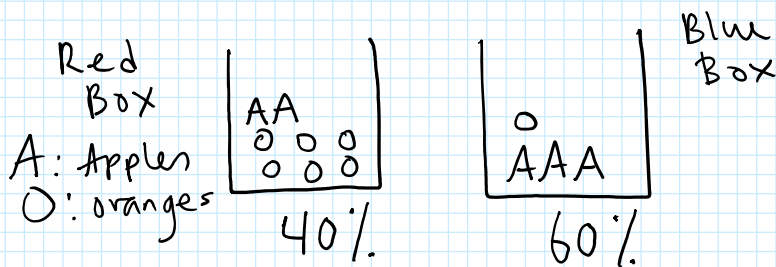
2- Select training parameters

3- train using training set

- 3 - train using training Set
 - 4 - Evaluate using Validation Set
 - 5 - Repeat steps 2, 3 & 4
 - 6 - Select the best Model & train with training & Validation Sets.
 - 7 - Estimate test error, $E^{(test)}$, using test set.
-

Data Normalization

- Standardization (z-score)
- for each of D dimensions
 - subtract mean
 - divide by standard deviation
- Use this mean & standard deviation to normalize data from test set.



$$P(B=r) = 0.4$$

$$P(B=b) = 0.6$$

$$P(F=A) = ? \quad P(B=r | F=O) = ?$$

B & F are Random Variables
discrete

$P(F)$ & $P(B)$ are prob. mass functions

Three Rules: Sum Rule Product Rule Bayes Rule

$$P(F=A | B=r) = \frac{1}{4}$$

$$P(F=O | B=r) = \frac{3}{4}$$

$$P(F=A | B=b) = \frac{3}{4}$$

$$P(F=O | B=b) = \frac{1}{4}$$

joint
prob.
mass
func.

$$\rightarrow P(F, B) = P(F|B) P(B) \leftarrow \text{product rule}$$

two events
occurring
simultaneously

$$\begin{aligned} P(F=A, B=r) &= P(F=A | B=r) \\ &\cdot P(B=r) \\ &= \frac{1}{4} \cdot \frac{4}{10} = \frac{1}{10} \end{aligned}$$

$$P(F) = \sum_B P(F, B)$$

2-D Histogram

$F \backslash B$	r	b
A	$\frac{1}{10}$	$\frac{9}{20}$
O	$\frac{3}{10}$	$\frac{3}{20}$

$P(F=A) = \frac{11}{20}$
 $P(F=O) = \frac{9}{20}$
 $P(B=r) = \frac{4}{10}$
 $P(B=b) = \frac{6}{10}$

Sum Rule

$$P(B) = \frac{1}{10} + \frac{1}{20} = \frac{3}{20}$$

$$P(F=A) = P(F=A, B=r) + P(F=A, B=b) = \frac{1}{10} + \frac{1}{20} = \frac{3}{20}$$

marginal prob.

likelihood prior marginalization

Posteriori

$$\rightarrow P(B|F) = \frac{P(F|B)P(B)}{P(F)} \leftarrow \text{Bayes Rule}$$

- Converts priors to posteriors
exploits role of evidence

$$P(B=r|F=0) = \frac{P(F=0|B=r)P(B=r)}{P(F=0)} = \frac{\frac{3}{4} \cdot \frac{1}{10}}{\frac{9}{20}} = \frac{2}{3}$$

Without evidence $P(B=r) = \frac{1}{10} = 0.1$

With " $P(B=r|F=0) = 0.67$

Statistical Independence between R.V. if

$$P(F, B) = P(F) \cdot P(B)$$

or $P(F|B) = P(F)$

for Continuous R.V.

Instead of Prob. Mass functions, we have
" density "

$p(x)$ prob. density func. for a R.V. X

$$\text{prob. } (x \in (a, b)) = \int_a^b p(x) dx$$

$$p(x) \geq 0$$

$$\int_{-\infty}^{\infty} p(x) = 1$$

Say we have two R.V. $X \neq Y$

Sum Rule $\rightarrow p(x) = \int p(x, y) dy$ Where $p(x, y)$ joint prob. density func.

product Rule $\rightarrow p(x, y) = p(x|y) p(y)$
 $= p(y|x) p(x)$

Expectation operator $\rightarrow E(\underline{x}) = \int x p(x) dx$ $E[f(\underline{x})] = \int f(x) p(x) dx$

Conditional expectation $E(x|y) = \int x p(x|y) dx \leftarrow \text{func. of } y$

$$p(x|y) = \frac{p(y|x) p(x)}{p(y)}$$

S.I. $\Leftrightarrow p(x, y) = p(x) \cdot p(y)$

$\begin{matrix} X \\ Y \end{matrix}$	129	130	131
15	0.12	0.42	0.06
16	0.08	0.28	0.04

$$p(X=130|Y=15) = ?$$

$$E(x) = ?$$

$$E(Y|X=129) = ?$$

Are $X \perp Y$ S.I.?

Given a set of samples $\{\underline{x}^{(1)}, \underline{x}^{(2)} \dots \underline{x}^{(N)}\}$
 we would like to estimate the joint prob. density (pdf)
 $P(\underline{x}^{(1)}, \underline{x}^{(2)}, \dots, \underline{x}^{(N)})$, assuming we know the form of pdf
 \rightarrow attempt to estimate its parameters, $\underline{\theta}$.

$\underline{x}^{(1)} \quad \underline{x}^{(2)} \quad \dots \quad \underline{x}^{(N)}$
 they came from a Normal distribution } Given
 What are μ & σ ?

$$\underline{\theta} = \begin{bmatrix} \mu \\ \sigma \end{bmatrix}$$

\rightarrow Assume \underline{x} 's are i.i.d.

$$P(\underline{x}^{(1)}, \underline{x}^{(2)}, \dots, \underline{x}^{(N)}; \underline{\theta}) = \prod_{i=1}^N P(\underline{x}^{(i)}; \underline{\theta})$$

this function is called the likelihood function
 of $\underline{\theta}$ with respect to examples & MLE finds
 $\underline{\theta}$ for which this function is maximum, i.e.,

$$\underline{\theta}_{ML} = \max_{\underline{\theta}} \prod_{i=1}^N P(\underline{x}^{(i)}; \underline{\theta})$$

equivalently

$$\underline{\theta}_{ML} = \max_{\underline{\theta}} \log \prod_{i=1}^N P(\underline{x}^{(i)}; \underline{\theta})$$

\log is a monotonically increasing func.
 $x_1 > x_2 \quad \log x_1 > \log x_2$

$$x_1 > x_2$$

$$\log x_1 > \log x_2$$

$$\underline{\theta}_{\text{ML}} = \max_{\underline{\theta}} \sum_{l=1}^N \underbrace{\log P(x^{(i)}; \underline{\theta})}_{\text{log-likelihood func.}}$$