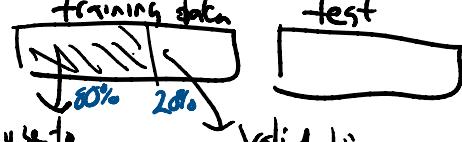


5.3-5.5: Hyperparameters and Validation Sets; Maximum Likelihood

Saturday, June 9, 2018 10:12 PM

hyperparameters that control model capacity cannot be learned in training
→ would always learn maximum possible capacity
 ↳ overfitting

validation set - data withheld from training
cannot use test set to decide hyperparameters



use to learn parameters validation guide selection of hyperparameters
estimate generalization error

benchmarks became stale

researchers built on results, so corrupted by use of valid, test data

(5.3.1) Cross-Validation

k-fold cross validation

for trial i, learn on data without subset i, test on subset i

(5.4) Estimators, Biases, Variance

(5.4.1) Point Estimation

point estimator (or statistic) is any function of the data:

$$\hat{\theta}_m = g(x^{(1)}, \dots, x^{(m)}) \quad \{x^{(1)}, \dots, x^{(m)}\}$$

set of i.i.d. points

goal is to approximate θ , true parameter value

function estimation

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

↑ not predictable from x

goal is to find \hat{f} that approximates f .

note: this is a type of point estimator, in function space

(5.4.2) Bias

- - - - - type of point estimator, in function space
(5.4.2) Bias

$$\text{bias}(\hat{\theta}_m) = E(\hat{\theta}_m) - \theta$$

expectation over ^{true} value
the data

estimator is **unbiased** if $\text{bias}(\hat{\theta}_m) = 0$

asymptotically unbiased if $\lim_{m \rightarrow \infty} \text{bias}(\hat{\theta}_m) = 0$.

$$\Rightarrow \lim_{m \rightarrow \infty} E(\hat{\theta}_m) = \theta.$$

(potential exercises, examples in book)

(5.4.3) Variance and Standard Error

$\text{Var}(\hat{\theta})$ standard error is $\sqrt{\text{Var}(\hat{\theta})} \approx SE(\hat{\theta})$

how much we expect the estimate to vary, as we resample from underlying data-generating process

std err of mean:

$$SE(\hat{\mu}_m) = \sqrt{\text{Var}\left[\frac{1}{m} \sum_{i=1}^m x^{(i)}\right]} = \frac{\sigma}{\sqrt{m}} \quad \text{where } \sigma^2 \text{ is true variance of } x^{(i)}$$

does not provide an unbiased estimate of standard deviation?

tends to underestimate it, but used in practice

often, in ML estimate generalization error as sample mean of error on the test set

central limit theorem: mean will be approx. distributed

e.g., 95% confidence interval around mean $\hat{\mu}_m$ is: with normal distribution

$$(\hat{\mu}_m - 1.96 SE(\hat{\mu}_m), \hat{\mu}_m + 1.96 SE(\hat{\mu}_m))$$

(5.4.4) Trade off Bias and Variance to Minimize MSE.

$$\text{MSE} = E[(\hat{\theta}_m - \theta)^2] = \text{Bias}(\hat{\theta}_m)^2 + \text{Var}(\hat{\theta}_m).$$

overall expected deviation
estimator true

overall expected estimator true deviation

increasing capacity tends to: decrease bias
increase variance

(5.4.5) Consistency

$$\text{plim}_{m \rightarrow \infty} \hat{\theta}_m = \theta \quad \text{consistency ("weak")}$$

converging in probability

$$P(|\hat{\theta}_m - \theta| > \epsilon) \rightarrow 0 \text{ as } m \rightarrow \infty$$

consistency ensures bias induced by estimator diminishes as m grows

(5.5) Maximum Likelihood Estimation

how can we derive good estimators?
functions that are

m examples drawn i.i.d. from $p_{\text{data}}(x)$

$$\mathbb{X} = \{x^{(1)}, \dots, x^{(m)}\}$$

$p_{\text{model}}(x; \theta)$ - parametric family of probability distributions
over space indexed by θ

so, $p_{\text{model}}(x; \theta)$ maps $x \rightarrow$ real number estimating probability $p_{\text{data}}(x)$.

maximum likelihood estimator for θ is defined as:

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta} p_{\text{model}}(\mathbb{X}; \theta) = \arg \max_{\theta} \prod_{i=1}^m p_{\text{model}}(x^{(i)}; \theta)$$

difficult to compute accurately (underflow)

so, take log: (doesn't change arg max):

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta} \sum_{i=1}^m \log p_{\text{model}}(x^{(i)}; \theta).$$

also, doesn't change when scaled \Rightarrow expectation

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta} E_{x \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(x; \theta). \quad (5.5g)$$

intuition: minimize dissimilarity between \hat{p}_{data} and p_{model} .

intuition: minimize dissimilarity between \hat{p}_{data} and p_{model} .
 empirical dist.
 defined by training set model distribution

KL divergence:

$$D_{\text{KL}}(\hat{p}_{\text{data}} \parallel p_{\text{model}}) = \mathbb{E}_{x \sim \hat{p}_{\text{data}}} [\log \hat{p}_{\text{data}}(x) - \log p_{\text{model}}(x)]$$

training a model to minimize KL divergence means minimize
 $-\mathbb{E}_{x \sim \hat{p}_{\text{data}}} [\log p_{\text{model}}(x)]$
 (same as maximizing S.Sq)

corresponds to minimizing cross-entropy between distributions

(5.5.1) Conditional Log-Likelihood and Mean Squared Error
 generalize Maximum Likelihood Estimation to estimate conditional probability, $P(y|x; \theta)$. (predict y given x)

X : all inputs, Y : all targets

Conditional MLE:

$$\Theta_{\text{ML}} = \arg \max_{\theta} P(Y|X; \theta)$$

assuming i.i.d.,

$$\Theta_{\text{ML}} = \arg \max_{\theta} \sum_{i=1}^m \log P(y^{(i)} | x^{(i)}; \theta)$$

Example: Linear Regression as Max Likelihood

explains why minimizing mean squared error

(5.5.2) Properties of Maximum Likelihood

main appeal: best estimator asymptotically as $m \rightarrow \infty$.
 rate of convergence

satisfies consistency: (converges to true value) under 2 conditions:

- P_{data} must lie within model family $p_{\text{model}}(\cdot; \theta)$
 (otherwise, no estimator can recover P_{data})

- p_{data} ... drawn from family $p_{\text{model}}(\cdot; \theta)$
(otherwise, no estimator can recover p_{data})
- true dist. of p_{data} must correspond to exactly one value of θ .
(otherwise, can recover, but can't know which θ is used)