

Statistics

LLN

$$\lim_{n \rightarrow \infty} \left| \frac{1}{n} \sum_{i=1}^n x_i - \mathbb{E}[x] \right| = 0$$

Unbiasedness

$$\text{Bias}_{\theta} = \mathbb{E} \left[\hat{\theta} \right] - \theta = \mathbb{E} \left[\hat{\theta} - \theta \right] = 0$$

Consistency

$$\lim_{n \rightarrow \infty} \mathbf{P} \left(\left| \hat{X} - \frac{1}{n} \sum_{i=1}^n X_i \right| \leq \epsilon \right) = 1$$

Properties of Gaussians

$$\begin{aligned} X &\sim \mathcal{N}(\mu, \sigma^2) \implies Y = aX + b \sim \mathcal{N}(a\mu + b, a^2\sigma^2) \\ X &\sim \mathcal{N}(\mu_x, \sigma_x^2), \quad Y \sim \mathcal{N}(\mu_y, \sigma_y^2), \implies Z = X + Y \sim \mathcal{N}(\mu_x + \mu_y, \sigma_x^2 + \sigma_y^2) \end{aligned}$$

Information Theory

Information

$$\mathbf{I}(S) = \log_2 \frac{1}{\mathbf{P}(S)}$$

Information Gain

$$\mathbf{I}(Y, X_i) = \mathbf{H}[Y] - \mathbf{H}[Y | X_i]$$

Entropy

$$\begin{aligned} \mathbf{H}[S] &= \sum_{s \in \{S\}} \mathbf{P}(S=s) \cdot \mathbf{I}(S=s) = \sum_{s \in \{S\}} \mathbf{P}(S=s) \log_2 \frac{1}{\mathbf{P}(S=s)} \\ &= \mathbb{E} \left[\log_2 \frac{1}{\mathbf{P}(S=s)} \right] = -\mathbb{E} \left[\log_2 \mathbf{P}(S=s) \right] \\ &= -\sum_{s \in \{S\}} \mathbf{P}(S=s) \log_2 \mathbf{P}(S=s) \end{aligned}$$

Learning Theory

PAC Learning

$$\mathbf{P} \left(\left| \hat{\theta} - \theta^* \right| \leq \epsilon \right) \geq 1 - \delta$$

$$\mathbf{P} \left(\left| \hat{\theta} - \theta^* \right| > \epsilon \right) < \delta$$

Decision Theory

Risk

Bayes Risk

$$\text{Risk}(f) = R(f) := \mathbb{E}_{(X,Y)} \left[\ell(Y, f(X)) \right]$$

$$R(f^*) \leq R(f), \quad \forall f \in \mathcal{F}$$

Bayes Optimal Rule

$$f^*(P) = \arg \min_f \mathbb{E}_{(X,Y) \sim P} \left[\ell(Y, f(X)) \right]$$

Emperical Risk Minimization

$$\hat{f}_n = \arg \min_f \frac{1}{n} \sum_{i=1}^n \left[\ell(Y_i, f(x_i)) \right] \xrightarrow[n \rightarrow \infty]{\text{LNN}} \arg \min_f \mathbb{E}_{(X,Y)} \left[\ell(Y, f(X)) \right]$$

Risk

True Risk

$$R(f) := \mathbb{E}_{(X,Y)} \left[\ell(f(X), Y) \right]$$

Emperical Risk

$$\hat{R}_{\mathcal{D}}(f) := \frac{1}{|\mathcal{D}|} \sum_{i \in \mathcal{D}} \left[\ell(f(X_i), Y_i) \right]$$

Excess Risk

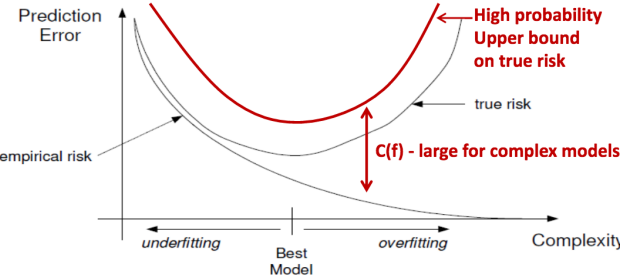
$$\mathbb{E} \left[R(\hat{f}_n) \right] - R(f^*)$$

Structural Risk

$$\left| R(f) - \hat{R}_n(f) \right| \leq C(f) \quad \forall f \in \mathcal{F}$$

$$R(f) \leq \hat{R}_n(f) + C(f) \quad \forall f \in \mathcal{F}$$

Use $\hat{R}_n(\hat{f}_n) + C(\hat{f}_n)$ as a *pes-simistic* estimate of true risk.



Risk Estimation

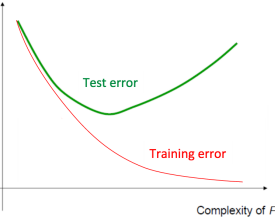
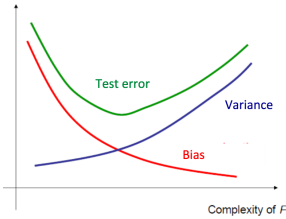
Bias–Variance Tradeoff

$$\text{Bias} = \mathbb{E} \left[f(X) \right] - f^*(X)$$

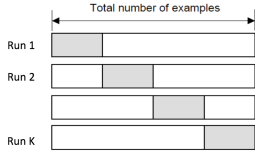
$$\text{Variance} = \mathbb{E} \left[(f(X) - \mathbb{E} \left[f(X) \right])^2 \right]$$

$$\text{Bias}^2 + \text{Variance} =$$

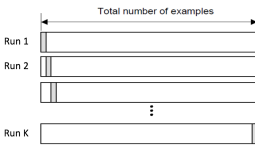
$$\mathbb{E} \left[(f(X) - f^*(X))^2 \right]$$



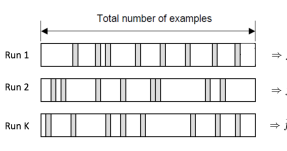
K–Fold CV



LOO CV



Random Subsampling



Hold–out Method

- split into two sets
- use \mathcal{D}_T to train a predictor $\hat{f}_{\mathcal{D}_T}$
- use \mathcal{D}_V to evaluate the predictor, $\hat{R}_{\mathcal{D}_V}(\hat{f}_{\mathcal{D}_T})$

$$\mathcal{D} = \{(X_i, Y_i)\}_{i=1}^n$$

↓

$$\mathcal{D}_T = \{(X_i, Y_i)\}_{i=1}^m$$

training set

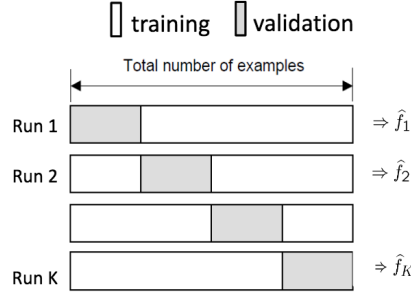
$$\mathcal{D}_V = \{(X_i, Y_i)\}_{i=m+1}^n$$

holdout set

Estimating True Risk

Estimate the error of a predictor on n data points. If K is large (close to n), bias of error estimate is small since each training set has close to n data points. However, variance of error estimates is high since each validation set has fewer data points and \hat{R}_{V_k} might deviate a lot from the mean.

$$\text{Error estimate} = \frac{1}{K} \sum_{k=1}^K \hat{R}_{V_k}(\hat{f}_{T_k})$$



Risk Minimization

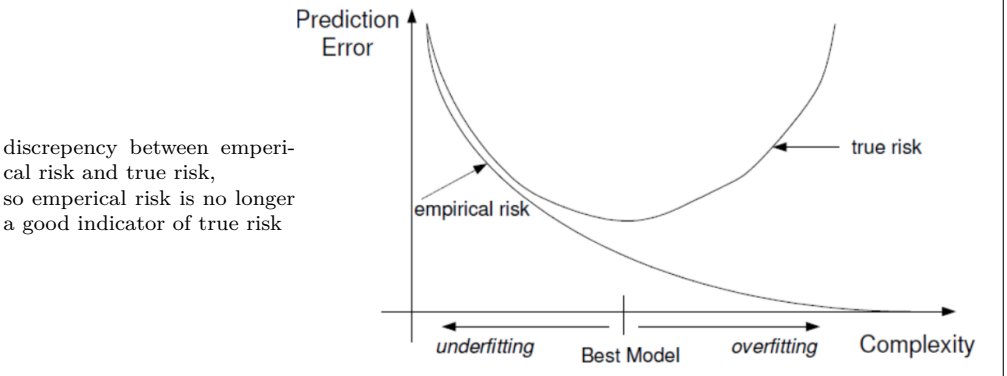
Emperical Risk Minimization

$$\hat{f}_n = \arg \min_f \frac{1}{n} \sum_{i=1}^n \left[\ell(Y_i, f(x_i)) \right]$$

Structural Risk Minimization

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}} \left[\hat{R}_n(f) + \lambda C(f) \right]$$

Overfitting



Regularization

Complexity Regularization

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}} \left\{ \hat{R}_n(f) + \lambda C(f) \right\}$$

$C(f)$: bound on deviation from true risk
 λ : regularization penalty

Information Criteria

AIC (Akiake IC) $C(f) = \# \text{ parameters}$

Allows # parameters to be infinite as # training data n becomes large

BIC (Bayesian IC) $C(f) = \# \text{ parameters} * \log n$

Penalizes complex models more heavily – limits complexity of models as # training data n becomes large

Model Selection

- define a finite set of model classes
- estimate true risk for each model class
- select model class with lowest estimated true risk

Model classes $\{\mathcal{F}_\lambda\}$ of increasing complexity $\mathcal{F}_1 < \mathcal{F}_2 < \dots$

$$\min_{\lambda} \min_{f \in \mathcal{F}_\lambda} J(f, \lambda)$$

- given λ estimate \hat{f}_λ using *empirical* / *structural* / *complexity regularized* risk minimization
- select λ for which \hat{f}_λ has minimum true risk estimated using *cross-validation* / *hold-out* / *information criteria*

Generalization Error

Terms

$$\text{estimated predictor} := \hat{f}_n$$

$$\text{optimal predictor} := f^*$$

$$\text{risk of estimated predictor} := R(\hat{f}_n)$$

$$\text{risk of optimal predictor} := R(f^*)$$

$$\text{expected risk} := \mathbb{E} \left[R(\hat{f}_n) \right]$$

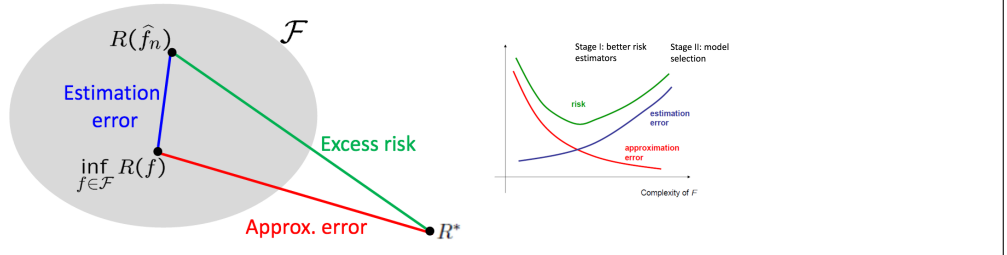
$$\text{excess risk} := \mathbb{E} \left[R(\hat{f}_n) \right] - R(f^*)$$

True Risk Decomposition

$$\mathbb{E} \left[R(\hat{f}_n) \right] - R^* = \underbrace{\left(\mathbb{E} \left[R(\hat{f}_n) \right] - \inf_{f \in \mathcal{F}} R(f) \right)}_{\text{estimation error}} + \underbrace{\left(\inf_{f \in \mathcal{F}} R(f) - R^* \right)}_{\text{approximation error}}$$

Estimation Error : due to randomness of training data

Approximation Error : due to restriction of model class



Regression
Linear Regression

Ridge Regression
$\hat{\theta}_{\text{MAP}} \arg \min_{\theta} \sum_{i=1}^n (Y_i - X_i \theta)^2 + \lambda \ \theta\ _2^2$

Lasso Regression
$\hat{\theta}_{\text{MAP}} \arg \min_{\theta} \sum_{i=1}^n (Y_i - X_i \theta)^2 + \lambda \ \theta\ _1$

Polynomial Regression

Classification
Logistic Regression

Naive Bayes

Boosting

Decision Trees
$\begin{aligned} \arg \max_{X_i} \left[\mathbf{H} \left[Y \right] - \mathbf{H} \left[Y \mid X_i \right] \right] &= \arg \min_{X_i} \mathbf{H} \left[Y \mid X_i \right] \\ &= \arg \min_{X_i} \sum_{x \in \{X_i\}} \left[\mathbf{P} \left(X_i = x \right) \mathbf{H} \left[Y \mid X_i = x \right] \right] \\ &= \arg \min_{X_i} - \sum_{x \in \{X_i\}} \left[\mathbf{P} \left(X_i = x \right) \sum_{y \in \{Y\}} \left[\mathbf{P} \left(Y = \mid X_i = x \right) \log_2 \mathbf{P} \left(Y = y \mid X_i = x \right) \right] \right] \end{aligned}$

Support Vector Machines

Deep Learning
Common Activation Functions

Backpropagation

Gradient Descent

Perceptron

MLP
CNN
RNN
LSTM
Clustering
KNN
Kernel Regression
Kernel Trick