

Machine Learning Cheat Sheet

General

Distributions

Multivariate gaussian distribution: $\mathcal{N}(\mathbf{X}|\mu, \Sigma)$

$$\implies p(\mathbf{X} = \mathbf{x}) =$$

$$(2\pi)^{-d/2} |\Sigma|^{-1/2} \exp\left[-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)\right]$$

Gaussian distribution: $\mathcal{N}(X|\mu, \sigma^2)$

$$\implies p(X = x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right)$$

Poisson distribution: $\mathcal{P}(X|\lambda)$

$$\implies p(X = k) = \frac{\lambda^k}{k!} \exp(-\lambda)$$

Some complexities

Product:	$O(NMD)$
Transpose:	$O(NM)$
Inversion:	$O(n^3)$

Some derivations

In respect to a vector z :

$$\begin{array}{l|l} AB & A \frac{\partial B}{\partial z} + \frac{\partial A}{\partial z} B \\ A^{-1} & -A^{-1} \frac{\partial A}{\partial z} A^{-1} \\ y^T A x & x^T A^T \frac{\partial y}{\partial z} + y^T A \frac{\partial x}{\partial z} \end{array}$$

In respect to a matrix A :

$$\begin{array}{l|l} |A| & |A| A^{-1} \\ x^T A x & X X^T \end{array}$$

Distribution properties

Let Y be any transformation and $\Sigma = \text{cov}(Y)$, then $\text{cov}(\alpha^T Y, \beta^T Y) = \beta^T \Sigma \alpha$.

If $X \sim \mathcal{N}(\mu, \Omega)$ and $Y = a + BX$, then

$Y \sim \mathcal{N}(a + B\mu, B\Omega B^T)$.

Gaussian conditioning

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mu, \Lambda^{-1})$$

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1})$$

$$\Downarrow$$

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mu + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\Lambda^{-1}\mathbf{A}^T)$$

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\Sigma\{\mathbf{A}^T\mathbf{L}(\mathbf{y} - \mathbf{b}) + \Lambda\mu\}, \Sigma)$$

$$\text{where } \Sigma = (\Lambda + \mathbf{A}^T\mathbf{L}\mathbf{A})^{-1}$$

Properties

Bayes rule: $p(A, B) = p(A|B)p(B) = p(B|A)p(A)$

A matrix \mathbf{M} is **positive semidefinite** $\iff \forall$ nonzero vector \mathbf{a} , we have $\mathbf{a}^T \mathbf{M} \mathbf{a} \geq 0$

If \mathbf{V} symmetric positive definite, then for all $\mathbf{P} \neq 0$, $\mathbf{P}^T \mathbf{V} \mathbf{P}$ is positive semi-definite (and even positive definite if \mathbf{P} is not singular).

Jensen's inequality applied to log:

$$\log(\mathbb{E}[X]) \geq \mathbb{E}[\log(X)]$$

$$\implies \log(\sum_x x \cdot p(x)) \geq \sum_x p(x) \log(x)$$

Matrix inversion lemma:

$$(\mathbf{PQ} + \mathbf{I}_N)^{-1} \mathbf{P} = \mathbf{P}(\mathbf{QP} + \mathbf{I}_M)^{-1}$$

$$\text{Useful derivative: } \frac{\partial \mathbf{x}^T \mathbf{B} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{B} + \mathbf{B}^T) \mathbf{x}$$

Optimization methods

Grid Search

Simply try values for all parameters at regular intervals. Complexity: $\mathcal{O}(M^D N D)$, where M is the number of values tried in each dimension.

Gradient Descent

$$\text{Update rule: } \beta^{(k+1)} = \beta^{(k)} - \alpha \frac{\partial \mathcal{L}(\beta^{(k)})}{\partial \beta}$$

Complexity: $\mathcal{O}(IND)$ where I is the number of iterations we take.

The gradient for MSE comes out as:

$$\frac{\partial \mathcal{L}}{\partial \beta} = -\frac{1}{N} \tilde{X}^T (\mathbf{y} - \tilde{X} \beta)$$

Newton's method

$$\text{General rule: } \beta^{(k+1)} = \beta^{(k)} - \alpha \mathbf{H}_{\mathbf{k}}^{-1} \frac{\partial \mathcal{L}(\beta^{(k)})}{\partial \beta}$$

where $\mathbf{H}_{\mathbf{k}}$ is the $D \times D$ Hessian at step k :

$$\mathbf{H}_{\mathbf{k}} = \frac{\partial^2 \mathcal{L}(\beta^{(k)})}{\partial \beta^2}$$

Complexity: $\mathcal{O}(IND^2 + ID^3)$, with the D^3 cost coming from the inversion of $\mathbf{H}_{\mathbf{k}}$.

For linear models, $\partial^2 \mathcal{L} / \partial x^2 = N^{-1} X^T X$.

Expectation-Maximization

$$\theta_{t+1} =$$

$$\arg \max_{\theta} \sum_{n=1}^N \mathbb{E}_{p(r_n | x_n, \theta^{(i)})} [\log p(x_n, r_n | \theta)].$$

Regression

Simple linear regression: $y_n \approx \beta_0 + \beta_1 x_{n1}$

Multiple linear regression:

$$y_n \approx f(\mathbf{x}_n) := \beta_0 + \beta_1 x_{n1} + \beta_2 x_{n2} + \dots + \beta_D x_{nD}$$

Linear basis function model

We can create more complex models while staying in the linear framework by transforming the inputs X of dimensionality D through a function $\phi: D \rightarrow M$.

$y_n = \beta_0 + \sum_{i=1}^M \beta_i \phi_i(\mathbf{x}_n) = \tilde{\phi}^T(\mathbf{x}_n^T) \beta$. The optimal β can be computed in closed form by $\beta = (\tilde{\Phi}^T \tilde{\Phi})^{-1} \tilde{\Phi}^T \mathbf{y}$ where $\tilde{\Phi}$ is a matrix with N rows and the n -th row is $[1, \phi_1(\mathbf{x}_n), \dots, \phi_M(\mathbf{x}_n)]$.

But note this requires $\tilde{\Phi}^T \tilde{\Phi}$ to be invertible (well-conditioned: $\tilde{\Phi}$ full column-rank).

$$\text{Ridge regression: } \beta_{\text{ridge}} = (\tilde{\Phi}^T \tilde{\Phi} + \lambda \mathbf{I})^{-1} \tilde{\Phi}^T \mathbf{y}$$

Cost functions

$$\text{Huber loss: } \mathcal{L}_{\delta}(a) = \begin{cases} \frac{1}{2} a^2 & \text{for } |a| \leq \delta, \\ \delta(|a| - \frac{1}{2} \delta), & \text{otherwise.} \end{cases}$$

Hinge Loss: $[t]_+ = \max(0, t) = \max_{\alpha \in [0, 1]} \alpha t$

Epsilon insensitive (for SVM classification):

$$\mathcal{L}_{\epsilon}(y, \hat{y}) = \begin{cases} 0 & \text{if } |y - \hat{y}| \leq \epsilon, \\ |y - \hat{y}| - \epsilon, & \text{otherwise.} \end{cases}$$

Kernel Ridge regression

$$\beta^* = (X^T X + \lambda I_D)^{-1} X^T \mathbf{y} = X^T (X X^T + \lambda I_N)^{-1} \mathbf{y} = X^T \alpha^*$$

We've

$$\beta^* = \arg \min_{\beta} \frac{1}{2} (\mathbf{y} - X \beta)^T (\mathbf{y} - X \beta) + \frac{\lambda}{2} \beta^T \beta$$

$$\alpha^* = \arg \max_{\alpha} -\frac{1}{2} \alpha^T (X X^T + \lambda I_M)^T \alpha + \alpha^T \mathbf{y}$$

Classification

$$\text{Logistic Function: } \sigma(t) = \frac{\exp(t)}{1 + \exp(t)}$$

$$\text{Derivative: } \frac{\partial \sigma(t)}{\partial t} = \sigma(t)[1 - \sigma(t)]$$

Classification with linear regression: Use $y = 0$ as class \mathcal{C}_1 and $y = 1$ as class \mathcal{C}_2 and then decide a newly estimated y belongs to \mathcal{C}_1 if $y < 0.5$.

Logistic Regression

$$\frac{\partial \mathcal{L}(\beta)}{\partial \beta} = -\tilde{X}^T [\sigma(\tilde{X} \beta) - \mathbf{y}]$$

$$\frac{\partial^2 \mathcal{L}}{\partial \beta^2} = X^T S X, \text{ with } S = \frac{\partial \sigma(X \beta)}{\partial X \beta} \text{ diagonal and}$$

$$S_{nn} = \sigma(\tilde{x}_n^T \beta) (1 - \sigma(\tilde{x}_n^T \beta))$$

There's no closed form, we can use gradient descent.

Generalized linear model

$$p(y|\eta) = \frac{h(y)}{Z} \exp[\eta^T \phi(y) - A(\eta)] \text{ with}$$

$$Z = \int h(y) \exp[\eta^T \phi(y) - A(\eta)] dy.$$

We've

- A link function g such that $E(\phi(y)) = \mu = g^{-1}(\eta)$.

$$\bullet E(\phi(\eta)) = \frac{\partial A(\eta)}{\partial \eta}$$

$$\bullet \text{var}(\phi(\eta)) = \frac{\partial^2 A(\eta)}{\partial \eta^2}$$

With $\eta_n = \tilde{x}_n^T \beta$, we've

$$\bullet \frac{\partial \mathcal{L}}{\partial \beta} = \tilde{X}^T (g^{-1}(\eta) - \phi(y))$$

$$\bullet \frac{\partial^2 \mathcal{L}}{\partial \beta^2} = X^T S X, \text{ where } S \text{ diagonal and } S_{nn} = \partial^2 A(\eta_n) / \partial \eta_n^2.$$

Cost functions

$$\text{RMSE: } \sqrt{\frac{1}{N} \sum_{n=1}^N [y_n - \hat{p}_n]^2}$$

Log-Loss:

$$-\frac{1}{N} \sum_{n=1}^N y_n \log(\hat{p}_n) + (1 - y_n) \log(1 - \hat{p}_n)$$

Probabilistic framework

$$\begin{array}{l|l} \text{Gaussian } p & \mathcal{L}_{lik} \hat{=} -\mathcal{L}_{MSE} \\ \text{Laplace } p & \mathcal{L}_{lik} \hat{=} -\mathcal{L}_{MAE} \end{array}$$

It is a sample approximation of the expected likelihood: $\mathcal{L}_{lik}(\beta) \approx E_y[p(y | \beta)]$

Bayesian networks

$$\underbrace{p(y, \theta)}_{\text{joint}} = \underbrace{p(y|\theta)}_{\text{likelihood}} \underbrace{p(\theta)}_{\text{prior}} = \underbrace{p(\theta|y)}_{\text{marg. likeli}} \underbrace{p(y)}_{\text{posterior}}$$

Distribution of the type $p_{\theta}(\mathbf{x}) = \prod^K p_{\theta}(x_k | Pa_k)$, where Pa_k is the set of parents of x_k .

$$\bullet m_{i \rightarrow a}(z_i) = p(z_i) \prod_{b \in \mathcal{N}(i) \setminus \{a\}} m_{b \rightarrow i}(z_i)$$

$$\bullet m_{a \rightarrow i}(z_i) = \sum_{\mathbf{z}_j: j \neq i} p(y_a | Pa_a = \mathbf{z}_j) \prod_{j \in \mathcal{N}(a) \setminus \{i\}} m_{j \rightarrow a}(\mathbf{z}_j)$$

$\sum_{\mathbf{z}_j: j \neq i}$ means the sumS over all possible neighbors values except those of z_i , and \mathbf{z}_j means the relevant values used for computation (thus take the product of messages from all neighbors except i and send the corresponding z_j value).

Kernel methods

$$(\mathbf{K})_{i,j} = \kappa(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j).$$

Linear	$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$
Polynomial	$\kappa(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^T \mathbf{x}_j + c)^d$
RBF	$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\ \mathbf{x}_i - \mathbf{x}_j\ ^2}{2\sigma^2}\right)$

Properties of a Kernel:

- \mathbf{K} should be symmetric: $\mathbf{K}^T = \mathbf{K}$
- \mathbf{K} should be positive semi-definite: \forall nonzero vector \mathbf{a} , $\mathbf{a}^T \mathbf{K} \mathbf{a} \geq 0$.

Neural Networks

A feed forward Neural Network is organized in K layers, each layer with $M^{(k)}$ hidden units $z_i^{(k)}$.

Activations $a_i^{(k)}$ are computed as the linear combination of the previous layer's terms, with weights $\beta^{(k)}$ (one $M^{(k-1)} \times 1$ vector of weights for each of the $M^{(k)}$ activations). Activations are then passed through a (possibly nonlinear) function h to compute the hidden unit $z_i^{(k)}$.

$$\mathbf{x}_n \xrightarrow{\beta_i^{(1)}} a_i^{(1)} \xrightarrow{h} z_i^{(1)} \xrightarrow{\beta^{(2)}} \dots \mathbf{z}^{(K)} = \mathbf{y}_n$$

Backpropagation

It's an algorithm which computes the gradient of the cost \mathcal{L} w.r.t. the parameters $\beta^{(k)}$.

Forward pass: compute a_i , z_i and \mathbf{y}_n from \mathbf{x}_n .

Backward pass: $\delta_n^{(k)} \triangleq \partial \mathcal{L} / \partial a_n^{(k)}$, compute $\delta_n^{(k-1)} = \text{diag}[\mathbf{h}'(\mathbf{a}_n^{(k)})] (\mathbf{B}^{(k)})^T \delta_n^{(k)}$

$$\text{Conclude: } \partial \mathcal{L} / \partial \mathbf{B}^{(k)} = \sum_n \delta_n^{(k)} \left(\mathbf{z}_n^{(k)} \right)^T.$$

Support Vector Machines

Search for the hyperplane separating the data such that the gap (margin) is biggest. It minimizes the following cost function ("hinge loss"):

$\mathcal{L}_{SVM}(\beta) = \sum_{n=1}^N [1 - y_n \phi_n \beta]_+ + \frac{\lambda}{2} \sum_{j=1}^M \beta_j^2$. This is convex but not differentiable. Using min-max theorem, we've:

$\max_{\alpha \in [0; C]^N} \alpha^T \mathbf{1} - \frac{1}{2} \alpha^T \mathbf{Y} \mathbf{K} \mathbf{Y} \alpha$, with constant $\alpha^T \mathbf{y} = 0$.

Min-max theorem: If $G(\alpha, \beta)$ is convex in β and concave in α , then $\min_{\beta} \max_{\alpha} G(\alpha, \beta) = \max_{\alpha} \min_{\beta} G(\alpha, \beta)$.

Clustering

K-Means

Minimize $\mathcal{L}(r, \mu) = \sum_{k=1}^K \sum_{n=1}^N r_{nk} \|x_n - \mu_k\|^2$. Given $\mu_k, r_n = \min_{j=1:K} \|x_n - \mu_j\|^2$ (use bizarre notation for r_{nk}).

Given $r_{nk}, \mu_k = \frac{\sum_{n=1}^N r_{nk} x_n}{\sum_{n=1}^N r_{nk}}$.

Gaussian Mixture Models

Assumptions: $p(r_n = k | \pi) = \pi_k$ (where $r_n \in 1 : K$, and $p(x_n | r_n = k, \mu, \Sigma) = \mathcal{N}(x_n | \mu_k, \Sigma_k) = \prod_{k=1}^K \mathcal{N}(x_n | \mu_k, \Sigma_k)^{r_{nk}}$

$p(\mathbf{x}_i | \theta) = \sum_{k=1}^K p(x_n, r_n = k | \theta) = \sum_{k=1}^K \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)$.

Taking the likelihood, we've

$\mathcal{L} = \sum_{n=1}^N \log \left(\sum_{k=1}^K \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k) \right)$.

Max using EM Take log lower bound with

$p_{kn} = \frac{\pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)}{\sum_{k=1}^K \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)}$ (expressed with old θ), we want the maximize

$\sum_{k=1}^K p_{kn} \log \left(\frac{\pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)}{p_{kn}} \right)$, where only p_{kn}

is expressed in the old θ . Therefore,

$$\mu_k^{(i+1)} = \frac{\sum_{n=1}^N p_{kn}^{(i)} x_n}{\sum_{n=1}^N p_{kn}^{(i)}}$$

$$\Sigma_k^{(i+1)} = \frac{\sum_{n=1}^N p_{kn}^{(i)} (x_n - \mu_k^{(i+1)})(x_n - \mu_k^{(i+1)})^T}{\sum_{n=1}^N p_{kn}^{(i)}}$$

$$\pi_k^{(i+1)} = \frac{1}{N} \sum_{n=1}^N p_{kn}^{(i)}$$

Matrix factorization

Minimize $L = \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K (x_{dn} - w_d^T z_n)^2 + \frac{\lambda_w}{2} \sum_{d=1}^D w_d^T w_d + \frac{\lambda_z}{2} \sum_{n=1}^N z_n^T z_n$. Easy, when the other is fixed. Thus

- $Z^T = (W^T W + \lambda_z I_M)^{-1} W^T X$
- $W^T = (Z^T Z + \lambda_w I_M)^{-1} Z^T X^T$

PCA

Use SVD: $X = USV^T$ (with U, V square matrices and orthonormal, S quasi-diagonal and non-negative).

We've $X v_i = \left(\sum_{j=1}^D s_j u_j v_j^T \right) v_i = s_j u_j$. We set $W = US^{1/2}$ and $Z = VS^{1/2}$.

Concepts

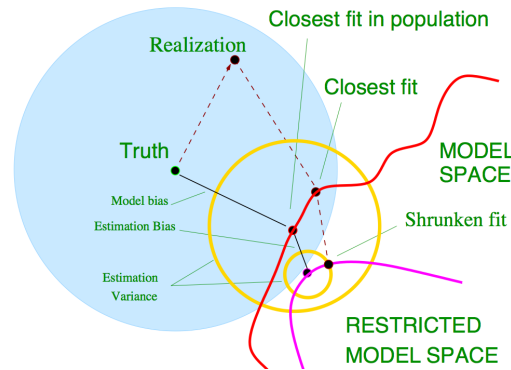
Convexity

f is called convex if: $\forall x_1, x_2 \in X, \forall t \in [0, 1] :$

$f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2)$.

Sum of two convex functions is convex. Composition of a convex function with a convex, nondecreasing function is convex. Linear, exponential and $\log(\sum \exp)$ functions are convex.

Bias-Variance Decomposition



Bias-variance comes directly out of the test error:

$$\begin{aligned} \overline{teErr} &= E[(\text{observation} - \text{prediction})^2] = E[(y - \hat{y})^2] \\ &= E[(y - y_{true} + y_{true} - \hat{y})^2] \\ &= E[(y - y_{true})^2] + E[(y_{true} - \hat{y})^2] \\ &\quad \text{var of measurement} \\ &= \sigma^2 + E[(y_{true} - E[\hat{y}] + E[\hat{y}] - \hat{y})^2] \\ &= \sigma^2 + E[(y_{true} - E[\hat{y}])^2] + E[(E[\hat{y}] - \hat{y})^2] \\ &\quad \text{pred bias}^2 \quad \text{pred variance} \end{aligned}$$

	bias	variance
regularization	+	-
choose simpler model	+	-
more data	-	

Identifiability

We say that a statistical model $\mathcal{P} = \{P_{\theta} : \theta \in \Theta\}$ is identifiable if the mapping $\theta \mapsto P_{\theta}$ is one-to-one:

$P_{\theta_1} = P_{\theta_2} \Rightarrow \theta_1 = \theta_2$ for all $\theta_1, \theta_2 \in \Theta$.

A non-identifiable model will typically have many local optima yielding the same cost, e.g.

$$\mathcal{L}(W, Z) = \mathcal{L}(aW, \frac{1}{a}Z)$$

Primal vs. Dual

Instead of working in the **column space** of our data, we can work in the **row space**:

$$\hat{\mathbf{y}} = \mathbf{X}\beta = \mathbf{X}\mathbf{X}^T \alpha = \mathbf{K}\alpha$$

where $\beta \in \mathbb{R}^D$ and $\alpha \in \mathbb{R}^N$ and (like magic) \mathbf{K} shows up, the Kernel Matrix.

Representer Theorem: For any β minimizing

$$\min_{\beta} \sum_{n=1}^N \mathcal{L}(y_n, \mathbf{x}_n^T \beta) + \sum_{d=1}^D \lambda \beta_d^2$$

there exists an α such that $\beta = \mathbf{X}^T \alpha$.

When we have an explicit vector formulation of β , we can use the matrix inversion lemma to get to the dual. E.g. for ridge regression:

$$\beta = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_D)^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{X}^T \underbrace{(\mathbf{X}\mathbf{X}^T + \lambda \mathbf{I}_N)^{-1}}_{\alpha} \mathbf{y}$$

In optimization, we get to the dual like this:

$$\begin{aligned} \min_{\beta} g(\beta) &\xrightarrow{(1)} \min_{\beta} \max_{\alpha} G(\beta, \alpha) \\ &\quad \downarrow (2) \\ \max_{\alpha} g^*(\alpha) &\xleftarrow{(3)} \max_{\alpha} \min_{\beta} G(\beta, \alpha) \\ &\quad \quad \quad g^*(\alpha) \end{aligned}$$

Gaussian process

Based on $(y, f_*)^T \sim \mathcal{N}\left(0, \begin{pmatrix} K & K_* \\ K_* & K_{**} \end{pmatrix}\right)$ then

$f_* | f \sim \mathcal{N}(\mu, \Sigma)$ s.t.

- $\mu = K_* K^{-1} y$
- $\Sigma = K_{**} - K_* K^{-1} K_*$

Random forest

We've $X = \{\mathbf{x}_i, y_i\}$. For each input variable k minimize $I_{\text{split}}(k)$: for $x_k > \tau$

- Split: $L = \{\mathbf{x}_i, y_i\}$ s.t. $x_{ik} > \tau$, and $R = X \setminus L$
- Compute probs: $p_L = \frac{\#(y_i=1)}{\#L}$ (for elem in L) and $p_R = \frac{\#(y_i=1)}{\#R}$ (for elem in R).
- $I_{\text{split}}(k) = \#L \cdot I(p_L) + \#R \cdot I(p_R)$, where I is an impurity measure

Different impurity measure:

- Misclassification error: $I(p) = 1 - \max(p, 1 - p)$
- Cross-entropy: $I(p) = -p \log(p) - (1 - p) \log(1 - p)$
- Gini-entropy: $I(p) = 2p(1 - p)$

Averaging if $\text{var}(f_i) = \sigma^2$ and $\text{corr}(f_i, f_j) = \rho \sigma^2$ (for $i \neq j$), with z_M the average of M trees, $\text{var}(z_M) = M^{-1} \sigma^2 + \rho \frac{M-1}{M} \sigma^2$.

Gain by small covariance

- Bagging: N observations with replication from X (don't change its size, just update some obs with replacement)
- Random. feature selection: subset of input variable (typi: $m_{\text{try}} = \sqrt{D}$)

Consistency

An estimator is said to be consistent, if it eventually recovers the true parameters that generated the data as the sample size goes to infinity. Of course, this only makes sense if the data actually comes from the specified model, which is not usually the case. Nevertheless, it can be a useful theoretical property.

Efficiency

An estimator is called efficient if it achieves the Cramer-Rao lower bound: $\text{Var}(\beta) \geq 1/I(\beta)$, where I is the Fisher information.

Occam's Razor

It states that among competing hypotheses, the one with the fewest assumptions should be selected. Other, more complicated solutions may ultimately prove correct, but in the absence of certainty, the fewer assumptions that are made, the better.

TODO: K-fold cross-validation, definition of Test-Error, Train-Error

TODO: statistical goodness (robust to outliers, ...) vs. computational goodness (convex, low computational complexity, ...) tradeoff. No free lunch theorem.

TODO: Decision Trees and Random Forests and Model averaging

Credits

Fork from Denwid's cheat-sheet.

Most material was taken from the lecture notes of Prof. Emteyaz Khan.

Cost functions figure from Patrick Breheny's slides.

Biais-variance decomposition figure from Hastie,

Tibshirani, Friedman, *The Elements of Statistical*

Learning. The SVD figure from Kevin P. Murphy,

Machine Learning, A Probabilistic Perspective.

Rendered January 12, 2016. Written by Dennis Meier and Merlin Nimier-David.
© Dennis Meier. This work is licensed under the Creative Commons Attribution-ShareAlike 3.0 Unported License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-sa/3.0/> or send a letter to Creative Commons, 444 Castro Street, Suite 900, Mountain View, California, 94041, USA.

