

Machine Learning Cheat Sheet

Regression

We have a set of N training examples of dimensionality D , e.g:

$$\mathbf{x}_n = [x_{n1} \ x_{n2} \ \dots \ x_{nD}]^T$$

We often put the examples into one matrix with extra column of 1s:

$$\tilde{\mathbf{X}} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1D} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N1} & x_{N2} & \dots & x_{ND} \end{bmatrix}$$

(note that \mathbf{x}_n are rows, not columns)

The goal is to predict a \hat{y} given \mathbf{x} .

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

Multiple dimension linear regression:

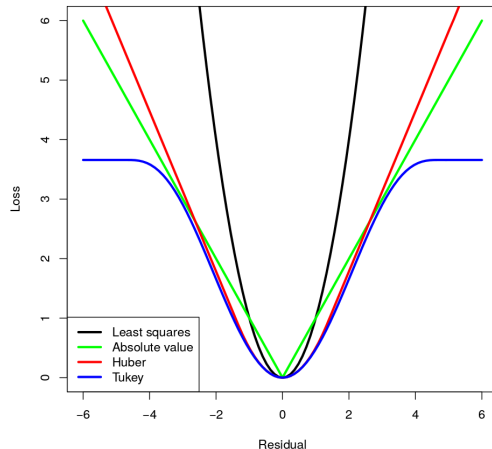
$$y_n \approx f(\mathbf{x}_n) := \beta_0 + \beta_1 x_{n1}^* + \beta_2 x_{n2}^* + \dots + \beta_D x_{nD}^*$$

Linear basis function model

$y_n = \beta_0 + \sum_{i=1}^M \beta_i \phi_i(\mathbf{x}_n) = \tilde{\phi}^T(\mathbf{x}_n) \beta$. The optimal β is given 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 $[1, \phi_1(x_n)^T, \dots, \phi_M(x_n)^T]$.

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

Cost functions



Cost function / Loss: $\mathcal{L}(\beta) = \mathcal{L}(\mathcal{D}, \beta)$

Mean square error (MSE): $\frac{1}{2N} \sum_{n=1}^N [y_n - f(\mathbf{x}_i)]^2$

Mean absolute error (MAE): $\frac{1}{2N} \sum_{n=1}^N |y_n - f(\mathbf{x}_i)|$

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

Root mean square error (RMSE): $\sqrt{2 * \text{MSE}}$

Epsilon insensitive (used for SVMs):

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

TODO: statistical/computational tradeoff

Gradient Descent

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

How to get a good α is a hard question (too small gives slow time, too big may not converge).

The gradient for MSE comes out as:

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

Normal equations

$$\beta = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{y}$$

Works for $\tilde{\mathbf{X}}$ having $\text{rank} = D$ (so that inverse is defined)

Classification

Logistic Function $\sigma = \frac{\exp(x)}{1 + \exp(x)}$

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

Logistic Regression

$$\tilde{\mathbf{X}}^T [\sigma(\tilde{\mathbf{X}} \beta) - \mathbf{y}] = 0$$

TODO: Generalized Linear model

Cost functions

Root Mean square error (RMSE):

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

0-1 Loss: $\frac{1}{N} \sum_{n=1}^N \delta(y_n, \hat{y}_n)$

logLoss:

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

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.

Math

convexity:

$$\forall x_1, x_2 \forall t \in [0, 1] f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2)$$

Jensen's inequality (log is concave):

$$\log(\sum_{i=1}^n x_i) \geq \sum_{i=1}^n \log(x_i)$$

Hessian is positive semidefinite \Rightarrow function is convex.

Positive semidefinite matrix $M \Leftrightarrow$ Gram matrix

(inner product) $\Leftrightarrow \forall x x^T M x \geq 0$.

Difficult words

consistent estimator converges to the true value as we increase the number of data to infinity.

unidentifiable model has many global minima due to symmetry.

Distributions

Gaussian: $\mathcal{N}(X|\mu, \sigma^2)$
 $\Rightarrow p(X = x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{1}{2}(\frac{x-\mu}{\sigma})^2)$

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

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

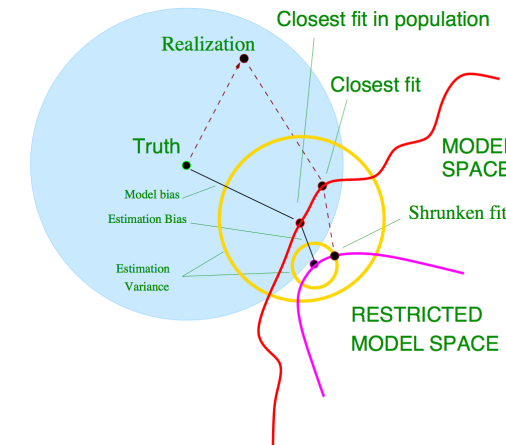
Complexities

Grid search: $\mathcal{O}(M^D ND)$, where M is the number of test points in one dimension.

Gradient descent: $\mathcal{O}(IND)$ where I is the number of iterations we make.

Least-squares (normal equations): $\mathcal{O}(ND^2 + D^3)$

Bias-Variance Decomposition



	bias	variance
regularization	+	-
reduce model complexity	+	-
more data	-	-

Maximum Likelihood

The Likelihood Function maps the model parameters to the probability distribution of \mathbf{y} :

$\mathcal{L}_{lik} : \text{parameter space} \rightarrow [0; 1] \quad \beta \mapsto p(\mathbf{y} | \beta)$ An underlying p is assumed before. If the observed \mathbf{y} are IID, $p(\mathbf{y} | \beta) = \prod_n p(\mathbf{y}_n | \beta)$.

\mathcal{L}_{lik} can be viewed as just another cost function. Maximum likelihood then simply chooses the parameters β such that observed data is most likely.

$$\beta = \arg \max_{\text{all } \beta} L(\beta)$$

Assuming different p is basically what makes this so flexible. We can choose e.g.:

Gaussian p	$\mathcal{L}_{lik} \hat{=} \mathcal{L}_{MSE}$
Poisson p	$\mathcal{L}_{lik} \hat{=} \mathcal{L}_{MAE}$

Bayesian methods

Bayes rule: $p(A, B) = p(A|B)p(B) = p(B|A)p(A)$
The **prior** $p(\mathbf{f}|\mathbf{X})$ encodes our prior belief about the "true" model \mathbf{f} . The **likelihood** $p(\mathbf{y}|\mathbf{f})$ measures the probability of our (possibly noisy) observations given the prior.

Least-squares tries to find model parameters β which maximize the likelihood. Ridge regression maximizes the **posterior** $p(\beta|\mathbf{y})$

Graphical Models

TODO: Bayes Net: Directed acyclic graph

TODO: Belief propagation

Graph between the observations and the variables is a bi-partite graph.

Kernel

Basically, Kernels are a mean to measure distance, or "similarity" of two vectors. We define:

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

The ϕ are not that important in the end, because we only use the Kernel as is. Sometimes it's even impossible to write them down explicitly.

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	$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\ \mathbf{x}_i - \mathbf{x}_j\ ^2}{2\sigma^2}\right)$

Neural Networks

TODO: Intuition and notation, backpropagation, regularization techniques

Support Vector Machines

Search for the hyperplane separating the data such that the gap is biggest. It minimizes the following cost function:

$$\mathcal{L}_{SVM}(\beta) = \sum_{n=1}^N [1 - y_n \tilde{\phi}_n \beta] + \frac{\lambda}{2} \sum_{j=1}^M \beta_j^2$$

This is convex but not differentiable.

Rendered January 5, 2015. Written by Dennis Meier, Jakub Sygnowski. © 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.

