

Machine Learning Cheat Sheet

TODO: statistical/computational tradeoff

Regression

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

$$x_n = [x_{n1} \ x_{n2} \ \dots \ x_{nD}]^T$$

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

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

(note that x_n are rows, not columns)

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

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

Multiple dimension linear regression:

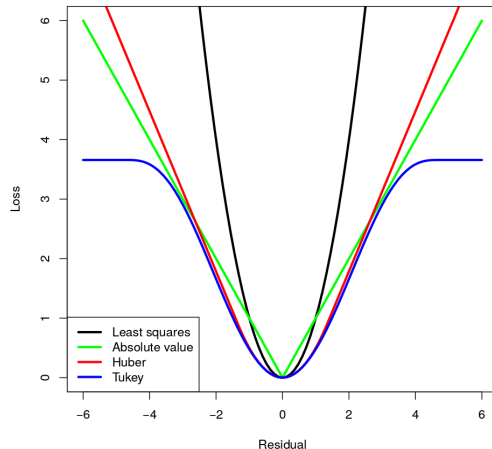
$$y_n \approx f(x^*) := \beta_0 + \beta_1 x_1^* + \beta_2 x_2^* + \dots + \beta_D x_D^*$$

Linear basis function model

$y_n = \beta_0 + \sum_{i=1}^M \beta_i \phi_i(\mathbf{x}_n) = \tilde{\Phi}^T(\mathbf{x}_n^T) \beta$. The optimal β is given by $\beta = (\tilde{\Phi}^T \tilde{\Phi})^{-1} \tilde{\Phi}^T y$ 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: $\beta_{\text{ridge}} = (\tilde{\Phi}^T \tilde{\Phi} + \lambda I)^{-1} \tilde{\Phi}^T y$

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):

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

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:

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

Newton's method

It uses second order derivatives information to converge faster (we approximate the objective function by a quadratic function rather than a line).

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

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

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

Normal equations

$$\frac{\partial \mathcal{L}}{\partial \beta} = -y^T X + X^T X \beta = 0$$

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

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

Classification

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

$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$

Logistic Regression

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

There's no closed form, because input is passed through a nonlinear function.

Cost functions

Root Mean square error (RMSE):

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

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

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):

$$\log\left(\frac{\sum_{i=1}^n x_i}{n}\right) \geq \frac{\sum_{i=1}^n \log(x_i)}{n} \text{ or } \log\left(\sum_x p(x)\right) \geq \sum_x p(x) \log(x)$$

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

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

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

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

Complexities

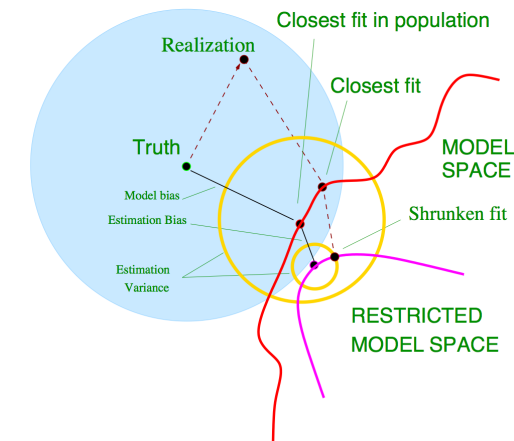
Grid search: $\mathcal{O}(M^D N D)$, where M is the number of test points in one dimension.

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

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

Newton's method (with Hessian): $\mathcal{O}(ND^2 + D^3)$

Bias-Variance Decomposition



Bias-variance come directly out of the test error:

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

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

Maximum Likelihood

The Likelihood Function maps the model parameters to the probability distribution of y :

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

\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. $\beta = \arg \max_{\beta} \mathcal{L}_{lik}(\beta)$

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	$\mathcal{L}_{lik} \hat{=} -\mathcal{L}_{MAE}$

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

TODO: put sample task

Mixture models

In mixture models, the data is generated by a sum (a mix) of K models. For GMM, these are gaussian:

$$\text{marginal: } p(x_i | \theta) = \sum_k p_k p(x_i, r_i = k | \theta) = \sum_k p_k p(x_i | r_i = k, \theta) \cdot \underbrace{p(r_i = k | \theta)}_{\pi_k}$$

$$\sum_{k=1}^K \pi_k p_k(x_i | \theta) = \sum_{k=1}^K \pi_k \mathcal{N}(x_i | \mu_k, \Sigma_k)$$

To use this for clustering, we first fit this mixture and then compute the posterior = $\frac{\text{likelihood} \times \text{prior}}{\text{marginal likelihood}}$

$$p_{nk}^{(c)} = p(r_i = k | x_i, \theta^{(c)}) = \frac{p(x_i | r_i = k, \theta^{(c)}) * p(r_i = k | \theta)}{\sum_k \pi_k p_k(x_i | \theta)}$$

max likelihood:

$$\max_{\theta} \log(p(x | \theta)) = \max_{\theta} \sum_n \log(p(x_n | \theta)) = \max_{\theta} \sum_n \log \sum_k \pi_k p_k(x_n | \theta)$$

We want to maximize (log of) the marginal:

e-step:

$$\log \sum_k \pi_k * p_k(x_n | \theta) = \log \sum_k \pi_k * p_k(x_n | \theta) / p_{kn}^{(c)}$$

$p_{kn}^{(c)} \geq (\text{Jensen}) \sum_k p_{kn}^{(c)} \log \pi_k * p_k(x_n|\theta)/p_{kn}^{(c)} = \log a + b - c (= \text{const})$
 m-step: $\max_{\theta} \sum_n \sum_k p_{kn}^{(c)} [\log \pi_k + \log p_k(x_n|\theta)]$ - maximization in regard of all parameters (π_k, μ_k, Σ_k) for GMM).

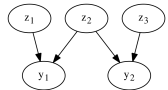
Bayesian methods

Bayes rule: $p(A, B) = p(A|B)p(B) = p(B|A)p(A)$
 The **prior** $p(f)$ 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.
 for Gaussian variable:

$$\begin{aligned} p(\mathbf{x}) &= \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \\ p(\mathbf{y}|\mathbf{x}) &= \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1}) \\ &\downarrow \\ p(\mathbf{y}) &= \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^T) \\ p(\mathbf{x}|\mathbf{y}) &= \mathcal{N}(\mathbf{x}|\boldsymbol{\Sigma}\{\mathbf{L}(\mathbf{y} - \mathbf{b}) + \mathbf{A}\boldsymbol{\mu}\}, \boldsymbol{\Sigma}) \\ \text{where } \boldsymbol{\Sigma} &= (\boldsymbol{\Lambda} + \mathbf{A}^T\mathbf{L}\mathbf{A})^{-1} \end{aligned}$$

Bayesian networks

We can use a Directed Acyclic Graph (DAG) to define a joint distribution of events. For example, we can express the relationship between *latent factors* (possible causes/outputs/variables) z_i and observations/features y_i :



This example can be factorized as follows:
 joint = $p(y_1, y_2, z_1, z_2, z_3) = p(y_1|z_1, z_2)p(y_2|z_2, z_3)p(z_1)p(z_2)p(z_3)$
 We can then obtain the z_i posterior distribution by marginalizing over the unknown variables:
 $p(z_1, z_2, z_3|y_1, y_2) = \frac{\text{joint}}{p(y_1, y_2)}$
 $\Rightarrow p(z_2|y_1, y_2) = \sum_{z_1, z_3} \underbrace{\frac{\text{joint}}{p(y_1, y_2)}}_{\text{message}} \propto \sum_{z_1, z_3} \text{joint}$

and to get marginals: $p(y) = \prod_i p(y_i)$
 $p(y_i) \stackrel{\text{e.g.}}{=} \sum_{z_2} p(y_i, z_2) = \sum_{z_2} p(y_i|z_2) \cdot p(z_2)$

Belief propagation

Belief propagation is a message-passing based algorithm used to compute desired marginals (e.g. $p(z_1|y_1, y_2)$) efficiently. It leverages the factorized expression of the joint. Messages passed:
 $m_{z_i \rightarrow y_j}(z_i) = p(z_i) \prod (\text{messages received by } z_i \text{ except from } y_j)$, e.g.
 $m_{z_1 \rightarrow y_2}(z_1) = p(z_1) m_{y_3 \rightarrow z_1}(z_1)$

F = set of all variables y_j depends upon
 $m_{y_j \rightarrow z_i}(z_i) = \sum_{F \neq z_i} p(y_j|F) \prod_{z \in F \setminus \{z_i\}} m_{z \rightarrow y_j}(z)$, e.g.
 $m_{y_2 \rightarrow z_1}(z_1) = \sum_{z_2} \sum_{z_4} p(y_2|z_1, z_2, z_4) m_{z_2 \rightarrow y_2}(z_2) m_{z_4 \rightarrow y_2}(z_4)$
 find a way: $p(z_i|y_1, y_2) \propto p(z_i) \cdot m_{y_1 \rightarrow z_i} \cdot m_{y_2 \rightarrow z_i}$ (if both exist). Then write messages from definition and write all paths from smallest messages to z_i .

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 (RBF).

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} , $\mathbf{a}^T \mathbf{K} \mathbf{a} \geq 0$.

Gaussian Process

The predicting function f is interpreted as a random variable with jointly gaussian prior: $\mathcal{N}(f|\mathbf{0}, \mathbf{K})$.
 Defining the Kernel matrix $\mathbf{K} = \kappa(\mathbf{X}, \mathbf{X})$ defines the prior. The key idea is, that if \mathbf{x}_i and \mathbf{x}_j are deemed by the kernel to be similar, then we expect the output of f at those points to be similar, too.
 We can sample functions from this random variable f and we can use prior + measurements to generate predictions.
 If we have measurements \mathbf{y} available, we get a joint distribution with the $\hat{\mathbf{y}}$ to be predicted:
 $\begin{bmatrix} \mathbf{y} \\ \hat{\mathbf{y}} \end{bmatrix} = \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} \kappa(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I} & \kappa(\mathbf{X}, \hat{\mathbf{X}}) \\ \kappa(\hat{\mathbf{X}}, \mathbf{X}) & \kappa(\hat{\mathbf{X}}, \hat{\mathbf{X}}) \end{bmatrix}\right)$
 This can be marginalized on \mathbf{y} to find the PDF of $\hat{\mathbf{y}}$.
 Advantage: we represent uncertainty.

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 $(M^{(k-1)} + 1) \times 1$ vector of weights for each of the $M^{(k)}$ activations - $\beta^{(k)}$ size is $M^{(k)} \times (M^{(k-1)} + 1)$ (one for free term, as in linear regression)). Activations are then passed through a (nonlinear) function h to compute the hidden unit $z_i^{(k)}$.
 $\mathbf{x}_n \xrightarrow{\beta_i^{(1)}} \mathbf{a}_i^{(1)} \xrightarrow{h} \mathbf{z}_i^{(1)} \xrightarrow{\beta^{(2)}} \dots \xrightarrow{h} \mathbf{z}^{(K)} = \mathbf{y}_n$

Backpropagation

Algorithm: Forward pass: compute a_i, z_i and \mathbf{y}_n from \mathbf{x}_n . Backward pass: work out derivatives of error to the target $\beta_i^{(k)}$ using the chain rule. Emti gives some complicated formulas for that, but regular calculations should work fine.
 $\delta^{(k-1)} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{(k-1)}} = \text{diag}[h'(\mathbf{a}^{(k-1)})] \mathbf{B}^{(k)T} \delta^{(k)}$
 $\frac{\partial \mathcal{L}}{\partial \mathbf{B}^{(1)}} = \delta^{(1)} \mathbf{x}^T$
 $\frac{\partial \mathcal{L}}{\partial \mathbf{B}^{(k)}} = \delta^{(k)} \mathbf{z}^{(k)T}$

Regularization

NN are not *identifiable* (existence of many local optima), so we initialize NN randomly.
 NN are universal density estimators (can fit any function), and thus prone to severe overfitting. Techniques used to reduce overfitting include early stopping (stop optimizing when test error starts increasing) and "weight decay" (i.e. L_2 regularization).

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"):
 $\mathcal{L}_{SVM}(\beta) = \sum_{n=1}^N [\mathbf{1} - \mathbf{y}_n \tilde{\phi}_n \beta] + \frac{\lambda}{2} \sum_{j=1}^M \beta_j^2$
 with $[t]_+ = \max(0, t)$
 This is convex but not differentiable. In the dual, the same problem can be formulated as:
 $\max_{\alpha \in [0, C]^N} \alpha^T \mathbf{1} - 1/2 \alpha^T Y K Y \alpha, \alpha^T \mathbf{y} = 0$

PCA

Find the eigenvectors of the covariance matrix $\mathbf{X}^T \mathbf{X}$ of the data. These form an orthonormal basis $\{\mathbf{w}_1, \dots, \mathbf{w}_N\}$ for the data in the directions that have highest variance. One can then use the first $L < D$ vectors to rebuild the data: $\hat{\mathbf{x}}_i = \mathbf{W} \mathbf{z}_i = \mathbf{W} \mathbf{W}^T \mathbf{x}_i$, with $\mathbf{W} = [\mathbf{w}_1; \dots; \mathbf{w}_L]$. This minimizes mean square error $\frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2$.

SVD

The same as with PCA, we can do with SVD:

$$\begin{array}{c} D \\ \boxed{} \\ N \end{array} = \begin{array}{c} D \quad N-D \\ \boxed{} \end{array} = \begin{array}{c} D \\ \boxed{} \end{array} = \begin{array}{c} D \\ \boxed{} \end{array}$$

$\mathbf{X} = \mathbf{U} \mathbf{S} \mathbf{V}^T$

The singular vals of a $N \times D$ matrix \mathbf{X} are the square roots of the eigenvalues of the $D \times D$ matrix $\mathbf{X}^T \mathbf{X}$

Curse of dimensionality

With the dimensionality increase, every data point becomes arbitrarily far from every other data point and therefore the choice of nearest neighbor becomes random. In high dimension, data only covers a tiny fraction of the input space, making generalization extremely difficult. Or in other words, the volume of the space increases so fast that the available data become sparse.

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:

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

In optimization, we get to the dual like this:

$$\begin{array}{ccc} \min_{\beta} g(\beta) & \xrightarrow{①} & \min_{\beta} \max_{\alpha} G(\beta, \alpha) \\ & & \downarrow ② \\ \max_{\alpha} g^*(\alpha) & \xleftarrow{③} & \max_{\alpha} \min_{\beta} G(\beta, \alpha) \\ & & \underbrace{\hspace{1cm}}_{g^*(\alpha)} \end{array}$$