

Computational Optimization

with Applications to Machine Learning

A Series of Lecture Notes at Missouri S&T

Wenqing Hu *

Contents

1	Motivating Examples	1
1.1	Supervised Learning	1
1.2	Matrix Optimizations	5

*Department of Mathematics and Statistics, Missouri University of Science and Technology (formerly University of Missouri, Rolla). Web: <https://huwenqing0606.github.io/> Email: huwen@mst.edu

1 Motivating Examples

1.1 Supervised Learning

EMPIRICAL RISK MINIMIZATION

Supervised Learning: Given training data points $(x_1, y_1), \dots, (x_n, y_n)$, construct a learning model $y = g(x, \omega)$ that best fits the training data. Here ω stands for the parameters of the learning model, say $\omega = (\omega_1, \dots, \omega_d)$.

Here (x_i, y_i) comes from an independent identically distributed family $(x_i, y_i) \sim p(x, y)$, where $p(x, y)$ is the joint density. The model $x \rightarrow y$ is a black-box $p(y|x)$, which is to be fit by $g(x, \omega)$.

“Loss function” $L(g(x, \omega), y)$, for example, can be $L(g(x, \omega), y) = (g(x, \omega) - y)^2$.

Empirical Risk Minimization (ERM):

$$\omega_n^* = \arg \min_{\omega} \frac{1}{n} \sum_{i=1}^n L(y_i, g(x_i, \omega)) . \quad (1.1)$$

Regularized Empirical Risk Minimization (R-ERM):

$$\omega_n^* = \arg \min_{\omega} \frac{1}{n} \sum_{i=1}^n L(y_i, g(x_i, \omega)) + \lambda R(\omega) . \quad (1.2)$$

For example, we can take $R(\omega) = \|\omega\|^2 = \omega_1^2 + \dots + \omega_d^2$. This regularization helps to control very irregular minimizers (unwanted ω).

In general let $f_i(\omega) = L(y_i, g(x_i, \omega))$ or $f_i(\omega) = L(y_i, g(x_i, \omega)) + \lambda R(\omega)$, then the optimization problem is

$$\omega_n^* = \arg \min_{\omega} \frac{1}{n} \sum_{i=1}^n f_i(\omega) . \quad (1.3)$$

Key features of nonlinear optimization problem in machine learning: large-scale, nonconvex, ... etc.

Key problems in machine learning: optimization combined with generalization.

“Population Loss”: $\mathbf{EL}(g(x, \omega), y)$, minimizer

$$\omega^* = \arg \min_{\omega} \mathbf{EL}(g(x, \omega), y) .$$

Generalization Error: $\mathbf{EL}(y, g(x, \omega_n^*))$. Consistency: Do we have $\omega_n^* \rightarrow \omega^*$? At what speed?

Key problems in optimization: convergence, acceleration, variance reduction.

How can optimization be related to generalization? There are quite abstract notions related to this topic, such as Vapnik–Chervonenkis dimension (VC dimension)

and Radmacher complexity, which we might touch later. Also, we want to look at the geometry of the loss landscape, which is closely related to neural network structure.

LOSS FUNCTIONS

Classification Problems: label $y = 1$ or -1 . Choice of Loss function $L(y, g)$, $y = 1, -1$. 0/1 Loss: $\ell_{0/1}(y, g) = 1$ if $yg < 0$ and $\ell_{0/1}(y, g) = 0$ otherwise.

(1) Hinge Loss.

$$L(y, g) = \max(0, 1 - yg) ; \quad (1.4)$$

(2) Exponential Loss.

$$L(y, g) = \exp(-yg) ; \quad (1.5)$$

(3) Cross Entropy Loss.

$$L(y, g) = - \left(I_{\{y=1\}} \ln \frac{e^g}{e^g + e^{-g}} + I_{\{y=-1\}} \ln \frac{e^{-g}}{e^g + e^{-g}} \right) . \quad (1.6)$$

This is to use $p(y) = \frac{e^{yg}}{e^{yg} + e^{-yg}}$, $y = \pm 1$ and the binary cross entropy as

$$-(I_{\{y=1\}} \ln p(1) + I_{\{y=-1\}} \ln p(-1)) .$$

Regression Problems: Choice of Loss function $L(y, g)$.

(1) L^2 -Loss.

$$L(y, g) = |y - g|_2^2 . \quad (1.7)$$

L^2 -norm: $|x|_2^2 = x_1^2 + \dots + x_d^2$

(2) L^1 -Loss.

$$L(y, g) = |y - g|_1 . \quad (1.8)$$

L^1 -norm: $|x|_1 = |x_1| + \dots + |x_d|$.

(3) L^0 -Loss.

$$L(y, g) = |y - g|_0 . \quad (1.9)$$

L^0 -norm: $|x|_0 = \#\{i : x_i \neq 0, 1 \leq i \leq d\}$.

Regularized (penalize) term $R(\omega)$:

L^1 -regularized $R(\omega) = |\omega|_1$;

L^0 -regularized $R(\omega) = |\omega|_0$.

LEARNING MODELS

(1) Linear regression: $g(x, \omega) = \omega^T x$. $g(x, \omega) = \frac{1}{1 + \exp(-\omega^T x)}$.

Least squares problem:

$$\min_{\omega \in \mathbb{R}^d} \frac{1}{2m} \sum_{j=1}^m (x_j^T \omega - y_j)^2 = \frac{1}{2m} \|A\omega - y\|_2^2 \quad (1.10)$$

Here training data $(x_1, y_1), \dots, (x_m, y_m)$ where $x_i \in \mathbb{R}^d$, $y \in \mathbb{R}$, and $A = \begin{pmatrix} x_1^T \\ \dots \\ x_m^T \end{pmatrix}$.

Tikhonov regularization:

$$\min_{\omega} \frac{1}{2m} \|A\omega - y\|_2^2 + \lambda \|\omega\|_2^2 . \quad (1.11)$$

LASSO (Least Absolute Shrinkage and Selection Operator):

$$\min_{\omega} \frac{1}{2m} \|A\omega - y\|_2^2 + \lambda \|\omega\|_1 . \quad (1.12)$$

See [4].

(2) Support Vector Machines (SVM):

Set-up: $x_j \in \mathbb{R}^n$, $y_j \in \{1, -1\}$. Separating hyperplane $\omega^T x + \beta = 0$ where $\omega \in \mathbb{R}^n$ and $\beta \in \mathbb{R}$.

Classification Problem: Goal is to find a hyperplane $\omega^T x + \beta = 0$ such that it classifies the two kinds of data points most efficiently. The signed distance of any point $x \in \mathbb{R}^n$ to the hyperplane is given by $r = \frac{\omega^T x + \beta}{\|\omega\|}$. If the classification is good enough, we expect to have $\omega^T x_j + \beta > 0$ when $y = 1$ and $\omega^T x_j + \beta < 0$ when $y = -1$. After rescaling ω and β , we can then formulate the problem as looking for optimal ω and β such that $\omega^T x_j + \beta \geq 1$ when $y_j = 1$ and $\omega^T x_j + \beta \leq -1$ when $y_i = -1$. The closest few data points that match these two inequalities are called “support vectors”. The distance to the separating hyperplane created by two support vectors of opposite type is

$$\left| \frac{1}{\|\omega\|} \right| + \left| \frac{-1}{\|\omega\|} \right| = \frac{2}{\|\omega\|} .$$

So we can formulate the following optimization problem

$$\max_{\omega \in \mathbb{R}^n, \beta \in \mathbb{R}} \frac{2}{\|\omega\|} \text{ such that } y_j(\omega^T x_j + \beta) \geq 1 \text{ for } j = 1, 2, \dots, m .$$

Or in other words we have the *constrained* optimization problem

$$\min_{\omega \in \mathbb{R}^n, \beta \in \mathbb{R}} \frac{1}{2} \|\omega\|^2 \text{ such that } y_j(\omega^T x_j + \beta) \geq 1 \text{ for } j = 1, 2, \dots, m . \quad (1.13)$$

“Soft margin”: We allow the SVM to make errors on some training data points but we want to minimize the error. In fact, we allow some training data to violate $y_j(\omega^T x_j + \beta) \geq 1$, so that ideally we minimize

$$\min_{\omega \in \mathbb{R}^n, \beta \in \mathbb{R}} \frac{1}{2} \|\omega\|^2 + C \sum_{j=1}^m \ell_{0/1}(y_j(\omega^T x_j + \beta) - 1) .$$

Here the 0/1 loss is $\ell_{0/1}(z) = 1$ if $z < 0$ and $\ell_{0/1}(z) = 0$ otherwise, and $C > 0$ is a penalization parameter. We can then turn the 0/1 loss to Hinge Loss, that is why Hinge Loss comes in:

$$\min_{\omega \in \mathbb{R}^n, \beta \in \mathbb{R}} \frac{1}{2} |\omega|^2 + C \sum_{j=1}^m \max(0, 1 - y_j(\omega^T x_j + \beta)) . \quad (1.14)$$

We can introduce “slack variables” $\xi_i \geq 0$ to introduce weights to classification errors in the above problem. This leads to “Soft margin SVM with slack variables”:

$$\min_{\omega \in \mathbb{R}^n, \beta \in \mathbb{R}} \frac{1}{2} |\omega|^2 + C \sum_{j=1}^m \xi_j \text{ s.t. } y_j(\omega^T x_j + b) \geq 1 - \xi_j, \xi_j \geq 0, j = 1, 2, \dots, m . \quad (1.15)$$

See [6].

(3) Neural Network: “activation function” σ .

Sigmoid:

$$\sigma(z) = \frac{1}{1 + \exp(-z)} ; \quad (1.16)$$

tanh:

$$\sigma(z) = \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)} ; \quad (1.17)$$

ReLU (Rectified Linear Unit):

$$\sigma(z) = \max(0, z) . \quad (1.18)$$

Vector-valued activation: if $z \in \mathbb{R}^n$ then $\sigma : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is defined by $(\sigma(z))_i = \sigma(z_i)$ where each $\sigma(z_i)$ is the scalar activation function.

Fully connected neural network prediction function

$$g(x, \omega) = a^T \left(\sigma \left(W^{(H)} (\sigma(W^{(H-1)} (\dots (\sigma(W^{(1)} x + b_1)) \dots) + b_{H-1}) + b_H) \right) \right) . \quad (1.19)$$

Optimization

$$\min_{\omega} \frac{1}{2} \sum_{i=1}^n (g(x_i, \omega) - y_i)^2 .$$

Many other different network structures that we do not expand here: convolutional, recurrent (Gate: GRU, LSTM), ResNet, Transformer, ...

Two layer (one hidden layer) fully connected ReLU neural network has specific loss

function structure: our $W^{(1)} = \begin{pmatrix} \omega_1^T \\ \dots \\ \omega_m^T \end{pmatrix}$ and

$$g(x, \omega) = \sum_{r=1}^m a_r \max(\omega_r^T x, 0) . \quad (1.20)$$

Optimization problem is given by

$$\min_{\omega} \frac{1}{2} \sum_{i=1}^n \left(\sum_{r=1}^m a_r \max(\omega_r^T x_i, 0) - y_i \right)^2 .$$

Non-convexity issues: see [2].

1.2 Matrix Optimizations

Many machine learning/statistical learning problems are related to matrix optimizations.

(1) Matrix Completion: Each A_j is $n \times p$ matrix, and we seek for another $n \times p$ matrix \hat{X} such that

$$\hat{X} = \arg \min_X \frac{1}{2m} \sum_{j=1}^m (\langle A_j, X \rangle - y_j)^2 \quad (1.21)$$

where $\langle A, B \rangle = \text{tr}(A^T B)$. We can think of the A_j as “probing” the unknown matrix X . In other words, we want the best X such that $\langle A_j, X \rangle \approx y_j$ for all $1 \leq j \leq m$.

(2) Nonnegative Matrix Factorization: If the full matrix $Y \in \mathbb{R}^{n \times p}$ is observed, then we seek for $L \in \mathbb{R}^{n \times r}$ and $R \in \mathbb{R}^{p \times r}$ such that

$$\min_{L, R} \|LR^T - Y\|_F^2 \text{ subject to } L \geq 0 \text{ and } R \geq 0 . \quad (1.22)$$

Here $\|A\|_F = (\sum \sum |a_{ij}|^2)^{1/2}$ is the Frobenius norm of a matrix A . This is used very often in recommendation systems (see [1]).

See [3] for an overview.

(3) Principle Component Analysis (PCA): Let S be a positive-definite (non-negative definite) matrix of size $n \times n$. Then we can diagonalize it as $Se_i = \lambda_i e_i$, $1 \leq i \leq n$, $\lambda_1 \geq \dots \geq \lambda_n > 0$ (or ≥ 0). Let $v \in \mathbb{R}^n$ be written as $v = v_1 e_1 + \dots + v_n e_n$. Then

$$v^T S v = \lambda_1 v_1^2 + \dots + \lambda_n v_n^2$$

is a quadratic form. If we restrict $v_1^2 + \dots + v_n^2 = 1$, then the maximum of above quadratic form will give the direction of e_1 and value λ_1 (principle component).

PCA:

$$\max_{v \in \mathbb{R}^n} v^T S v \text{ such that } \|v\|_2 = 1, \|v\|_0 \leq k . \quad (1.23)$$

Here S is a positive-definite (or non-negative definite) matrix. The objective function is convex, but if you take into account the constraint, then this problem becomes non-convex. A picture for dimension 1 example can be shown below.

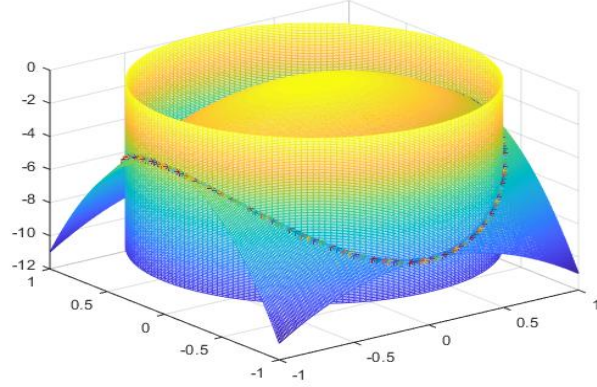


Figure 1: Loss Landscape of 1-dimensional PCA.

Online PCA: see [5].

(4) Sparse inverse covariance matrix estimation: Sample covariance matrix $S = \frac{1}{m-1} \sum_{j=1}^m a_j a_j^T$. $S^{-1} = X$. “Graphical LASSO”:

$$\min_{X \in \text{Symmetric } \mathbb{R}^{n \times n}, X \succeq 0} \langle S, X \rangle - \ln \det X + \lambda \|X\|_1 \quad (1.24)$$

where $\|X\|_1 = \sum |X_{ij}|$.

References

- [1] [https://en.wikipedia.org/wiki/Matrix_factorization_\(recommender_systems\)](https://en.wikipedia.org/wiki/Matrix_factorization_(recommender_systems)).
- [2] Alon Brutzkus, Amir Globerson, Eran Malach, and Shai Shalev-Shwartz. Sgd learns over-parameterized networks that provably generalize on linearly separable data. *ICLR, arXiv:1710.10174*, 12:1–12, 2018.
- [3] Y. Chi, Y.M. Lu, and Y. Chen. Nonconvex optimization meets low-rank matrix factorization: An overview. *IEEE Transactions on Signal Processing*, 67(20), 2019.
- [4] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, 2009.
- [5] C.J Li, M. Wang, H. Liu, and T. Zhang. Near-optimal stochastic approximation for online principal component estimation. *Mathematical Programming*, 167(1):75–97, 2018.
- [6] Vladimir N. Vapnik. *The Nature of Statistical Learning Theory*. Springer, 2000.