

Matrix Methods in Machine Learning

Lecture Notes

Rebekah Dix

November 11, 2018

Contents

| | | |
|----------|---|-----------|
| 1 | Elements of Machine Learning | 3 |
| 2 | Linear Algebra Review | 3 |
| 2.1 | Products | 3 |
| 2.2 | Linear Independence | 4 |
| 3 | Linear Systems and Vector Norms | 5 |
| 3.1 | Vector Norms | 5 |
| 4 | Least Squares | 6 |
| 4.1 | Geometric Approach | 7 |
| 4.2 | Vector Calculus Approach | 8 |
| 4.2.1 | Review of Vector Calculus | 8 |
| 4.2.2 | Application to Least Squares | 9 |
| 4.3 | Positive Definite Matrices | 9 |
| 4.4 | Subspaces | 10 |
| 4.5 | Least Squares with Orthonormal Basis for Subspace | 10 |
| 4.5.1 | Orthogonal Matrices and Orthonormal Basis | 10 |
| 4.5.2 | Back to LS | 11 |
| 4.5.3 | Gram-Schmidt Orthogonalization Algorithm | 11 |
| 5 | Least Squares Classification | 13 |
| 6 | Tikhonov Regularization/Ridge Regression | 14 |
| 6.1 | Tikhonov Regularization Derivation | 15 |
| 6.1.1 | Derivation with Vector Calculus | 15 |
| 6.1.2 | Alternative Derivation | 15 |

| | | |
|-----------|---|-----------|
| 7 | Singular Value Decomposition | 16 |
| 7.1 | Interpretation of SVD | 17 |
| 7.2 | Low-Rank Approximation | 17 |
| 8 | Power Iteration and Page Rank | 18 |
| 8.1 | SVD: Connection to Eigenvalues/vectors | 18 |
| 9 | Matrix Completion | 18 |
| 9.1 | Iterative Singular Value Thresholding | 18 |
| 10 | Iterative Solvers | 18 |
| 10.1 | Gradient Descent/Landweber Iteration | 19 |
| 11 | Regularized Regression | 21 |
| 11.1 | Proximal Gradient Algorithm | 21 |
| 11.2 | LASSO (Least absolute selection and shrinkage operator) | 21 |
| 12 | Convexity and Support Vector Machines | 23 |
| 12.1 | Convexity | 23 |
| 12.2 | Support Vector Machines | 24 |

1 Elements of Machine Learning

1. Collect data
2. Preprocessing: changing data to simplify subsequent operations without losing relevant information.
3. Feature extraction: reduce raw data by extracting features or properties relevant to the model.
4. Generate training samples: a large collection of examples we can use to learn the model.
5. Loss function: To learn the model, we choose a loss function (i.e. a measure of how well a model fits the data)
6. Learn the model: Search over a collection of candidate models or model parameters to find one that minimizes the loss on training data.
7. Characterize generalization error (the error of our predictions on new data that was not used for training).

2 Linear Algebra Review

2.1 Products

Inner products:

$$\langle x, w \rangle = \sum_{j=1}^p w_j x_j = x^T w = w^T x \quad (1)$$

Thus this inner product is a weighted sum of the elements of x .

Matrix-vector multiplication:

$$Xw = \begin{bmatrix} -x_1^T- \\ -x_2^T- \\ \vdots \\ -x_n^T- \end{bmatrix} w = \begin{bmatrix} x_1^T w \\ x_2^T w \\ \vdots \\ x_n^T w \end{bmatrix} \quad (2)$$

Matrix-matrix multiplication:

Example 1. Let $X \in \mathbb{R}^{n \times p}$, n movies, p people. $T \in \mathbb{R}^{n \times r}$, and $W \in \mathbb{R}^{r \times p}$. We can think of T as the taste profiles of r representative customers and W as the weights on each representative profile (there will be one set of weights for each customer). Suppose we have two representative taste profiles (i.e. an action lover and a romance lover). Then w will be a 2-vector containing the weights of on the two representative taste profiles. Then

Tw is the expected preferences of a customer who weights the representative taste profiles of T with the weights given in w .

Now we can think about the full matrix product $X = TW$

$$X = TW \implies X_{ij} = \langle \text{ith row of } T, j\text{th column of } W \rangle \quad (3)$$

- The j th column of X is a weighted sum of the columns of T , where the j th column of W tells us the weights.

$$x_j = Tw_j \quad (4)$$

That is, the tastes (preferences) of the j th customer.

- The i th row of X is $x_i^T = t_i^T W$ where t_i^T is the i th row of T . This gives us how much each customer likes movie i .

Inner product representation:

$$TW = \begin{bmatrix} -t_1^T- \\ -t_2^T- \\ \vdots \\ -t_n^T- \end{bmatrix} \begin{bmatrix} | & | & \dots & | \\ w_1 & w_2 & \dots & w_p \\ | & | & & | \end{bmatrix} = \begin{bmatrix} t_1^T w_1 & t_1^T w_2 & \dots & t_1^T w_p \\ t_2^T w_1 & \ddots & & \vdots \\ \vdots & & \ddots & \vdots \\ t_n^T w_1 & & & t_n^T w_p \end{bmatrix} \quad (5)$$

Outer Product Representation:

$$TW = \begin{bmatrix} | & | & \dots & | \\ T_1 & T_2 & \dots & T_r \\ | & | & & | \end{bmatrix} \begin{bmatrix} -w_1^T- \\ -w_2^T- \\ \vdots \\ -w_r^T- \end{bmatrix} = \sum_{k=1}^r T_k w_k^T \quad (6)$$

(the sum of rank 1 matrices. TW has rank r if and only if the columns of T are rows of W are linearly independent). In this representation, we can think about T_k as the k th representative taste profile and w_k^T as the k th row of W , or the affinity of each customer with the k th representative profile.

2.2 Linear Independence

Definition 1. (Linear Independence) Vectors $v_1, v_2, \dots, v_n \in \mathbb{R}^p$ are linearly independent vectors if and only if

$$\sum_{j=1}^n \alpha_j v_j = 0 \iff \alpha_j = 0, j = 1, \dots, n \quad (7)$$

Definition 2. (Matrix rank) The rank of a matrix is the maximum number of linearly independent columns. The rank of a matrix is less than the smallest dimension of the matrix.

3 Linear Systems and Vector Norms

Example 2. (Condition on $\text{rank}(A)$ for existence of exact solution)

Consider the linear system of equations $Ax = b$. This means that b is a weighted sum of the columns of A . Suppose A is full rank. Now consider the matrix $[A \ b]$. If the rank of $[A \ b]$ were *greater* than the rank of A (since the number of columns of the matrix increased by 1 and A is assumed full rank, this would imply the rank is $\text{rank}(A) + 1$), this would mean that b could not be written as a linear combination of the columns of A , and that the system would not have an exact solution. Therefore, we must have that $\text{rank}([A \ b]) = \text{rank}(A)$ in order for the system $Ax = b$ to have an exact solution.

To see how the definition of linear independence applies here, observe that $Ax = b \implies Ax - b = 0$. Therefore

$$[A \ b] \begin{bmatrix} x \\ -1 \end{bmatrix} = 0 \quad (8)$$

Thus, if $Ax = b$ has an exact solution, then $[A \ b]$ does not have linearly independent columns.

Example 3. (Condition on $\text{rank}(A)$ for more than one exact solution)

If the system of linear equations $Ax = b$ has more than one exact solution, then there is at least one non zero vector w for which $x + w$ is also a solution. That is, $A(x + w) = b$. If x is an exact solution, then $Ax = b$. This implies $Aw = 0$. Therefore, the columns of A are linearly dependent. Thus, if $\text{rank}(A) < \dim(x)$, then there will be more than one exact solution.

Example 4. (Apply the above conditions) Let

$$A = \begin{bmatrix} 1 & -2 \\ -1 & 2 \\ -2 & 4 \end{bmatrix}, \quad b = \begin{bmatrix} 2 \\ -2 \\ -4 \end{bmatrix} \quad (9)$$

We want to solve $Ax = b$.

- This system has an exact solution, since $\text{rank}(A) = \text{rank}([A \ b])$. This follows since the columns of A are linearly dependent, so it has rank 1, and b is a multiple of the columns of A , so the rank of $[A \ b]$ is also 1.
- Note that $1 = \text{rank}(A) < \dim(x) = 2$. Therefore this system does not have a unique solution.

3.1 Vector Norms

Definition 3. (Vector Norm) A vector norm is a function $\|\cdot\|$ mapping from $\mathbb{R}^n \rightarrow \mathbb{R}$ with the following properties.

1. $\|x\| \geq 0$ for all $x \in \mathbb{R}^n$.

2. $\|x\| = 0$ if and only if $x = 0$.
3. $\|\alpha x\| = |\alpha| \|x\|$ for all $\alpha \in \mathbb{R}, x \in \mathbb{R}^n$.
4. $\|x + y\| \leq \|x\| + \|y\|$ for all $x, y \in \mathbb{R}^n$.

Helpful fact: $\|x\|_{q'} \leq \|x\|_q$ if $1 \leq q \leq q' \leq \infty$.

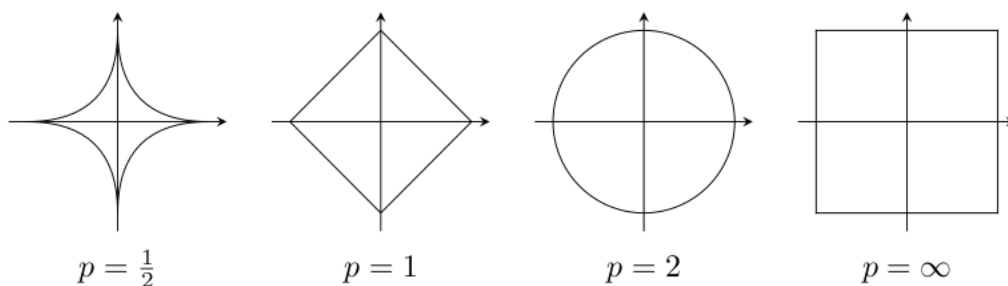


Figure 1: The l_p norm in \mathbb{R}^2

4 Least Squares

We are given:

1. Vector of labels $y \in \mathbb{R}^n$
2. Matrix of features $X \in \mathbb{R}^{n \times p}$

We want to find:

1. Vector of weights $w \in \mathbb{R}^p$

Assumptions:

1. $n \geq p$, and $\text{rank}(X) = p$.

If $y = Xw$, then we have a system of n linear equations, where the i th equation is

$$y_i = w_1 x_{i1} + w_2 x_{i2} + \cdots + w_p x_{ip} = \sum_{j=1}^p w_j x_{ij} = \langle w, x_{\cdot i} \rangle \quad (10)$$

where $x_{\cdot i}$ is the i th row of X .

In general, $y \neq Xw$ for any w . We define a residual $r_i = y_i - \langle w, x_{\cdot i} \rangle$. Our goal is then to find w $\sum_{i=1}^n |r_i|^2$ (the sum of square residuals/errors).

Why should we minimize the sum of square errors?

1. Magnifies the effect of large errors

2. Allows us to compute derivatives
3. Simple geometric interpretation
4. Coincides with modeling $y = Xw + \epsilon$, where ϵ is Gaussian noise

4.1 Geometric Approach

We know $\hat{r} = y - X\hat{w}$ is orthogonal to the span of the columns of X . Thus $x_i^T \hat{r} = 0$, or $X^T \hat{r} = 0$. This implies $X^T(y - X\hat{w}) = 0$. Thus \hat{w} is a solution to the linear system of equations

$$X^T X \hat{w} = X^T y \quad (11)$$

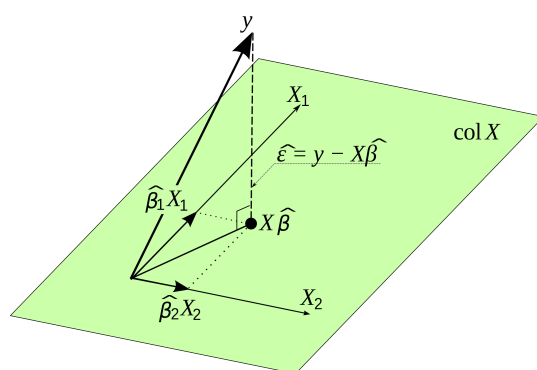


Figure 2: Geometry of LS in \mathbb{R}^2

Observations:

- The question we're trying to answer: What is the point in $col(X)$ that has the shortest distance to y ? In \mathbb{R}^2 , what are the weights β_1 and β_2 such that $\beta_1 x_1 + \beta_2 x_2$ has the shortest distance to y ?
- $col X$ is the space of all vectors that can be written as $\alpha x_1 + \beta x_2$ for some $\alpha, \beta \in \mathbb{R}$, that is the span of the columns of X . y may not lie in this space.
- The residual vector will form a right angle with $col X$, because any other angle would correspond to a longer distance.

4.2 Vector Calculus Approach

4.2.1 Review of Vector Calculus

Let w be a p -vector and let f be a function of w that maps \mathbb{R}^p to \mathbb{R} . Then the gradient of f with respect to w is

$$\nabla_w f(w) = \begin{pmatrix} \frac{\partial f(w)}{\partial w_1} \\ \vdots \\ \frac{\partial f(w)}{\partial w_p} \end{pmatrix} \quad (12)$$

Example 5. (Gradient of an Inner Product) Let $f(w) = \langle a, w \rangle = w^T a = \sum_{i=1}^n w_i a_i$. Then

$$\nabla_w w^T a = \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_p \end{pmatrix} = a \quad (13)$$

Example 6. (Gradient of an Inner Product, Squared) Let $f(w) = \|w\|^2 = w^T w = w_1^2 + \dots + w_p^2$. Then

$$\nabla_w w^T w = \begin{pmatrix} 2w_1 \\ 2w_2 \\ \vdots \\ 2w_p \end{pmatrix} = 2w \quad (14)$$

(This is a special case of the Quadratic Form discussed below, where $w^T Q w$, and $Q = I$)

Example 7. (Gradient of a Quadratic Form) Let $x \in \mathbb{R}^n$ and $f(x) = x^T Q x$, where Q is symmetric (if Q isn't symmetric we could replace Q with $\frac{1}{2}(Q + Q^T)$). Then

$$\begin{aligned} f(x) &= x^T Q x \\ &= \sum_{i=1}^n \sum_{j=1}^n x_i Q_{ij} x_j \end{aligned}$$

Therefore

$$[\nabla_x f]_k = \frac{df}{dx_k} = \begin{cases} 2Q_{kk}x_k & i = j = k \\ Q_{kj}x_j & i = k, i \neq j \\ Q_{ik}x_i & j = k, j \neq i \end{cases} \quad (15)$$

Therefore

$$\nabla_x f = (Q + Q^T)x \quad (16)$$

If Q is symmetric, then this equals $2Qx$.

4.2.2 Application to Least Squares

Let $f(w) = \|y - Xw\|_2^2$. Then the least squares problem is

$$\hat{w} = \arg \min_w f(w) \quad (17)$$

We can expand $f(w)$ as

$$\begin{aligned} f(w) &= (y - Xw)^T (y - Xw) \\ &= y^T y - y^T Xw - w^T X^T y + w^T X^T Xw \\ &= y^T y - 2w^T X^T y + w^T X^T Xw \end{aligned}$$

Then

$$\nabla_w f(w) = -2X^T y + 2X^T Xw$$

At an optimum we have that \hat{w} solves $X^T y = X^T X\hat{w}$. Then if $(X^T X)^{-1}$ exists, we have that

$$\hat{w} = (X^T X)^{-1} X^T y \quad (18)$$

Theorem 1. (*Sufficient Condition for Existence/Uniqueness of LS Solution*) If the columns of X are linearly independent, then $X^T X$ is non-singular, and there exists a unique least squares solution $\hat{w} = (X^T X)^{-1} X^T y$.

Proof. □

4.3 Positive Definite Matrices

Definition 4 (Positive Definite, pd). A matrix Q ($n \times n$) is positive definite (written $Q \succ 0$) if $x^T Q x > 0$ for all $x \in \mathbb{R}^n$, $x \neq 0$.

Definition 5 (Positive Semi-Definite, psd). A matrix Q ($n \times n$) is positive semi-definite (written $Q \succeq 0$) if $x^T Q x \geq 0$ for all $x \in \mathbb{R}^n$, $x \neq 0$.

Properties of Positive Definite matrices:

1. If $P \succ 0$ and $Q \succ 0$, then $P + Q \succ 0$.
2. If $P \succ 0$ and $\alpha > 0$, then $\alpha P \succ 0$.
3. For any matrix A , $A^T A \succeq 0$ and $AA^T \succeq 0$. Further, if the columns of A are linearly independent, then $A^T A \succ 0$.
4. If $A \succ 0$, then A^{-1} exists.
5. Notation: $A \succ B$ means $A - B \succ 0$.

Example 8. Let

$$X = \begin{pmatrix} 1 & 1 \\ 1 & 1 \\ 1 & 1 \end{pmatrix} \quad (19)$$

Then

$$X^T X = \begin{pmatrix} 3 & 3 \\ 3 & 3 \end{pmatrix} \quad (20)$$

Consider the vector $a = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$. Then $a^T X^T X a = 0$. Therefore $X^T X$ is not positive definite.

4.4 Subspaces

Definition 6. (Subspace) A set of points $S \subseteq \mathbb{R}^n$ is a subspace if

1. $0 \in S$ (S contains the origin)
2. If $x, y \in S$, then $x + y \in S$
3. If $x \in S, \alpha \in \mathbb{R}$, then $\alpha x \in S$.

4.5 Least Squares with Orthonormal Basis for Subspace

Suppose are given a training sample $\{x_i, y_i\}_{i=1}^n$, $x_i \in \mathbb{R}^p$ and $y \in \mathbb{R}$. If the columns of X (the data matrix) are linearly dependent, then $X^T X$ is not invertible. It is then impossible to tell which features are significant predictors of y .

Given $X = [x_1 \dots x_p]$, the following are options to represent the corresponding subspace spanned by the columns of X :

1. Use X , but then LS can be hard to interpret.
2. Use an orthonormal basis for the subspace.

4.5.1 Orthogonal Matrices and Orthonormal Basis

Definition 7. (Orthonormal basis for X) An orthonormal basis for the columns of X is a collection of vectors $\{u_1, \dots, u_r\}$ such that the span of the columns of X equals the span of $\{u_1, \dots, u_r\}$. That is, $\text{span}(\{x_1, \dots, x_p\}) = \text{span}(\{u_1, \dots, u_r\})$. Furthermore,

$$u_i^T u_j = \begin{cases} 0 & i \neq j \\ 1 & i = j \end{cases} \quad (21)$$

That is, the u vectors are orthogonal and have norm 1.

Observations:

- The rank r of the subspace must satisfy $r \leq \min(n, p)$. r is the number of linearly independent columns of X .
- We can place the basis vectors into a basis matrix $U \in \mathbb{R}^{n \times r}$.

Claim 1. (Properties of orthogonal (basis) matrices) Let $U \in \mathbb{R}^{n \times r}$ be an orthogonal (basis) matrix.

1. $U^T U = I$
2. If U and V are both orthogonal, then UV is also orthogonal.
3. U is length preserving: $\|Uv\|_2 = \|v\|_2$ for $v \in \mathbb{R}^n$.

Proof. We prove each item as follows:

1. We can easily see this from the inner product interpretation of matrix multiplication.
2. $(UV)^T UV = V^T U^T UV = V^T V = I$.
3. $\|Uv\|_2^2 = (Uv)^T Uv = v^T U^T Uv = v^T v = \|v\|_2^2$.

□

4.5.2 Back to LS

Suppose U is an orthonormal basis matrix for our data matrix X . Then, the least-squares problem is

$$\hat{v} = \arg \min_v \|y - Uv\|_2^2 \quad (22)$$

We need \hat{v} to satisfy $U^T U \hat{v} = U^T y$. Thus, $\hat{v} = U^T y$.

4.5.3 Gram-Schmidt Orthogonalization Algorithm

How can we take X and get an orthonormal basis U ?

1. Input $X = [x_1 \dots x_p] \in \mathbb{R}^{n \times p}$
Output: $U = [u_1 \dots u_r] \in \mathbb{R}^{n \times r}$
where $r = \text{rank}(X) \leq \min(n, p)$
2. Initialize $u_1 = \frac{x_1}{\|x_1\|_2}$
3. For $j = 2, 3, \dots, p$
 $x'_j =$ all the components of x_j not represented by u_1, \dots, u_{j-1} .

$$x'_j = x_j - \sum_{i=1}^{j-1} (u_i^T x_j) u_i \quad (23)$$

here $(u_i^T x_j)$ is the least squares weight for u_i .

$$u_j = \begin{cases} \frac{x'_j}{\|x'_j\|_2} & x'_j \neq 0 \\ 0 & x'_j = 0 \end{cases} \quad (24)$$

Next, by construction, each column of U , u_i , is in $\text{span}(\{x_1, \dots, x_p\})$. Therefore we can write

$$u_i = \alpha_{i1}x_1 + \alpha_{i2}x_2 + \dots + \alpha_{ip}x_p \quad (25)$$

where the $\alpha_{ij} \in \mathbb{R}$. We can write this in matrix form as

$$U = XA \quad (26)$$

where X is $n \times p$ and A is $p \times r$, and the i th column of A is

$$a_i = \begin{bmatrix} \alpha_{i1} \\ \alpha_{i2} \\ \vdots \\ \alpha_{ip} \end{bmatrix} \quad (27)$$

Thus, $u_i = Xa_i$.

Now, suppose $w \in \mathbb{R}^p$ is the vector of weights we found using LS, and as above, v is our vector of weights founding using LS with an orthonormal basis matrix. We have two equations for the predicated label \hat{y}

$$\begin{aligned} \hat{y} &= w_1x_1 + w_2x_2 + \dots + w_px_p \\ &= v_1u_1 + v_2u_2 + \dots + v_ru_r \\ &= v_1Xa_1 + v_2Xa_2 + \dots + v_rXa_r \\ &= v_1(\alpha_{11}x_1 + \alpha_{12}x_2 + \dots + \alpha_{1p}x_p) \\ &\quad \vdots \\ &\quad + v_r(\alpha_{r1}x_1 + \alpha_{r2}x_2 + \dots + \alpha_{rp}x_p) \\ &= x_1(v_1\alpha_{11} + \dots + v_r\alpha_{r1}) \\ &\quad \vdots \\ &\quad + x_p(v_1\alpha_{1p} + \dots + v_r\alpha_{rp}) \end{aligned}$$

Notice that

$$\begin{aligned} w_1 &= v_1\alpha_{11} + \dots + v_r\alpha_{r1} \\ &\quad \vdots \\ w_p &= v_1\alpha_{1p} + \dots + v_r\alpha_{rp} \end{aligned}$$

Therefore

$$\hat{y} = XAv = Xw \quad (28)$$

so that $Av = w$.

In sum, given a new sample $x_{new} \in \mathbb{R}^p$, we have two ways to predict label y_{new} :

1. $\hat{y}_{new} = \langle x_{new}, w \rangle$
2. Using an orthonormal basis U , we know that $U = XA$. Therefore, $u_{new}^T = x_{new}^T A$. Equivalently, $u_{new} = Ax_{new}$. Then $y_{new} = \langle u_{new}, v \rangle$.

If the columns of X are linearly independent ($r = p$), we can calculate using LS (recalling $u_i = Xa_i$)

$$a_i = (X^T X)^{-1} X^T u_i \quad (29)$$

Theorem 2. Let $X \in \mathbb{R}^{n \times p}$, $n \geq p$, be full rank (the p columns of X are linearly independent) and $y \in \mathbb{R}^n$. Let u_1, \dots, u_p be orthonormal basis vectors such that $\text{span}(\{x_1, \dots, x_p\}) = \text{span}(\{u_1, \dots, u_p\})$. Then $\hat{y} = X\hat{w}$ where $\hat{w} = \arg \min_w \|y - Xw\|_2^2$ is given by $\hat{y} = UU^T y$, where $U = [u_1 \ u_2 \ \dots \ u_p]$.

Proof.

$$\hat{y} = X\hat{w} = X(X^T X)^{-1} X^T y \quad (30)$$

where $P_x = X(X^T X)^{-1} X^T$ is a projection matrix. Since $\text{span}(\{x_1, \dots, x_p\}) = \text{span}(\{u_1, \dots, u_p\})$, we must have that

$$P_x y = P_u y \quad (31)$$

which implies $P_x = P_u$. Thus

$$P_x = P_u = U(U^T U)^{-1} U^T = UU^T \quad (32)$$

Finally

$$\hat{y} = P_x y = P_u y = UU^T y \quad (33)$$

□

5 Least Squares Classification

We are given a training sample $\{x_i, y_i\}_{i=1}^n$, $x_i \in \mathbb{R}^p$ and $y \in \mathbb{R}$ (or $y \in \{+1, -1\}$).

Definition 8. (Linear Predictor) We have a linear predictor if each label is a linear combination of the features i.e. we can find weights $\{w_i\}_{i=1}^p$ such that

$$y_i = w_1 x_{i1} + w_2 x_{i2} + \dots + w_p x_{ip} \quad (34)$$

In words, this says the label for observation i is a linear combination of the features for example i .

The steps to complete least squares classification in this environment are as follows:

1. Build a data matrix or feature matrix and label vector

$$X = \begin{bmatrix} -x_1^T - \\ -x_2^T - \\ \vdots \\ -x_n^T - \end{bmatrix} = \begin{bmatrix} x_1^T & 1 \\ x_2^T & 1 \\ \vdots & \vdots \\ x_n^T & 1 \end{bmatrix} \in \mathbb{R}^{n \times p}, \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad (35)$$

The linear model is then $\hat{y} = Xw$.

2. Solve a least squares optimization problem

$$\hat{w} = \arg \min_w \|y - Xw\|_2^2 = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2 \quad (36)$$

(this last equality makes it clear that we are minimizing the sum of squared residuals). If the columns of X are linearly independent, then $X^T X$ is positive definite. Therefore $X^T X$ is invertible. In sum, if $X^T X$ is positive definite, then there exists a unique LS solution

$$\hat{w} = (X^T X)^{-1} X^T y \quad (37)$$

The predicted labels are

$$\begin{aligned} \hat{y} &= Xw \\ &= X(X^T X)^{-1} X^T y \end{aligned}$$

3. Validate with test/hold out data

6 Tikhonov Regularization/Ridge Regression

We are given $X \in \mathbb{R}^{n \times p}$ (n training samples, p features) and $y \in \mathbb{R}^n$ (n labels). Our model is $y \approx Xw$, which means $y_i \approx x_i^T w$ for some $w \in \mathbb{R}^p$.

The LS problem is

$$\hat{w}_{LS} = \arg \min_w \|y - Xw\|_2^2 = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2 \quad (38)$$

There are two cases

1. If X is full rank (i.e. the columns of X are linearly independent), then \hat{w}_{LS} is unique and

$$\hat{w}_{LS} = (X^T X)^{-1} X^T y \quad (39)$$

2. If X is not full rank, then $X^T X$ is not invertible. \hat{w}_{LS} is not unique; there are infinitely many solutions.

6.1 Tikhonov Regularization Derivation

In this second case (and it can also be useful in the first), we can define a new objective

$$\hat{w} = \arg \min_w \|y - Xw\|_2^2 + \lambda \|w\|_2^2 \quad (40)$$

where $\|y - Xw\|_2^2$ measures the fit to the data, $\lambda > 0$ is a regularization parameter or tuning parameter, and $\|w\|_2^2$ is a regularizer. $\|w\|_2^2$ measures the energy in w .

Observations about this problem:

1. \hat{w} is unique even when no unique least square solution exists
2. Even when X is full rank, $X^T X$ can be badly behaved, and regularization adjusts for this.

6.1.1 Derivation with Vector Calculus

Let $f(w) = \|y - Xw\|_2^2 + \lambda \|w\|_2^2$. Then

$$\begin{aligned} f(w) &= y^T y - 2w^T X^T y + w^T X^T X w + \lambda w^T w \\ &= y^T y - 2w^T X^T y + w^T (X^T X + \lambda I) w \end{aligned}$$

Then

$$\nabla_w f(w) = -2X^T y + 2(X^T X + \lambda I) \quad (41)$$

If $(X^T X + \lambda I)$ is invertible, then $\hat{w} = (X^T X + \lambda I)^{-1} X^T y$. BUT, $(X^T X + \lambda I)$ is *always* invertible. Recall that if a matrix is positive definite, then it is invertible. We can show that $(X^T X + \lambda I)$ is indeed positive definite and hence invertible. To see this, fix $0 \neq a \in \mathbb{R}^n$, then

$$\begin{aligned} a^T (X^T X + \lambda I) a &= a^T X^T X a + \lambda a^T a \\ &= \|Xa\|_2^2 + \lambda \|a\|_2^2 \end{aligned}$$

Now note that $\|Xa\|_2^2 \geq 0$ (it could be 0 if X is not full rank and a is in the null space of X – this is what causes troubles with LS) but $\lambda \|a\|_2^2 > 0$. Therefore, $(X^T X + \lambda I)$ is positive definite.

6.1.2 Alternative Derivation

Note that for vectors a, b ,

$$\|a\|_2^2 + \|b\|_2^2 = \left\| \begin{bmatrix} a \\ b \end{bmatrix} \right\|_2^2 \quad (42)$$

Therefore,

$$\begin{aligned}
f(w) &= \|y - Xw\|_2^2 + \lambda \|w\|_2^2 \\
&= \|y - Xw\|_2^2 + \|\sqrt{\lambda}w\|_2^2 \\
&= \left\| \begin{bmatrix} y - Xw \\ \sqrt{\lambda}w \end{bmatrix} \right\|_2^2 \\
&= \left\| \begin{bmatrix} y \\ 0 \end{bmatrix} - \begin{bmatrix} Xw \\ \sqrt{\lambda}w \end{bmatrix} \right\|_2^2 \\
&= \left\| \begin{bmatrix} y \\ 0 \end{bmatrix} - \begin{bmatrix} X \\ \sqrt{\lambda}I \end{bmatrix} w \right\|_2^2 \\
&= \|\tilde{y} - \tilde{X}w\|_2^2
\end{aligned}$$

We can solve this problem with LS, so that

$$\hat{w} = (\tilde{X}^T \tilde{X})^{-1} \tilde{X} \tilde{y} \quad (43)$$

where

$$\tilde{X}^T \tilde{X} = X^T X + \lambda I \quad (44)$$

and

$$\tilde{X} \tilde{y} = X^T y \quad (45)$$

Thus this is equivalent to the derivation above.

7 Singular Value Decomposition

Theorem 3. Every matrix $X \in \mathbb{R}^{n \times p}$ can be factorized as

$$X = U \Sigma V^T \quad (46)$$

where

1. $U \in \mathbb{R}^{n \times n}$ is orthogonal ($U^T U = U U^T = I$) are the left singular vectors of X
2. $\Sigma \in \mathbb{R}^{n \times p}$ is diagonal and contains the singular values of X
3. $V \in \mathbb{R}^{p \times p}$ is orthogonal ($V^T V = V V^T = I$) and contains the right singular vectors of X

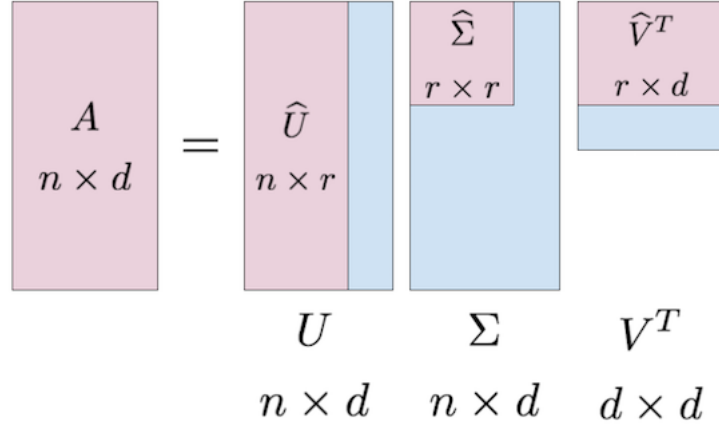


Figure 3: SVD

7.1 Interpretation of SVD

1. U is an orthonormal basis for the columns of X
2. ΣV^T are the basis coefficients

Example 9. (Netflix) Let $X \in \mathbb{R}^{n \times p}$ be a matrix (full rank) where the columns are taste profiles of customers and the rows are single movie ratings across customers.

1. The i th column of U is a basis vector in \mathbb{R}^n and is the i th representative customer taste profile (i.e. vector of normalized movie ratings).
2. The j th column of V^T (the j th row of V) is the relative importance of each representative taste profile to predicting customer j 's preferences.
3. The i th row of V^T (the i th column of V) is the vector of users' affinities to the i th representative profile.

7.2 Low-Rank Approximation

Theorem 4. (Subspace Approximation) If $X \in \mathbb{R}^{p \times n}$ has rank $r > k$, then

$$\min_{Z : \text{rank}(Z)=k} \|X - Z\|_F^2 \quad (47)$$

is given by $Z = X_k = U_k \Sigma_k V_k^T$ and

$$\|X - X_k\|_F^2 = \sum_{i=k+1}^r \sigma_i^2 \quad (48)$$

8 Power Iteration and Page Rank

8.1 SVD: Connection to Eigenvalues/vectors

Suppose $X = U\Sigma V^T \in \mathbb{R}^{p \times n}$. Then

$$\begin{aligned} A &:= X^T X \\ &= V\Sigma U^T U\Sigma V^T \\ &= V\Sigma^2 V^T \\ &= V\Lambda V^T \end{aligned}$$

Thus V is a matrix of the eigenvectors of A and $\Lambda = \Sigma^2$ contains its eigenvalues.

Power iteration gives a method to find the 1st right singular vector.

9 Matrix Completion

9.1 Iterative Singular Value Thresholding

Algorithm 1 Iterative Singular Value Thresholding

Require: $\hat{X} = \text{zeros}(n, p)$

Require: $\hat{X}_\Omega = X_\Omega$

▷ fill in obs. entries

Require: Threshold or r

for $k \leftarrow 0, 1, \dots$ **do**

$[U, S, V] \leftarrow \text{svd}(\hat{X})$

$\hat{S} \leftarrow S \geq \text{threshold}$

▷ if threshold, keep sing vals \geq

$\hat{S} \leftarrow S(1:r, 1:r)$

▷ if rank, keep r sing vals

$\hat{X} \leftarrow U\hat{S}V^T$

$\hat{X}_\Omega \leftarrow X_\Omega$

▷ fill in obs. entries

 If converged: $\|\hat{X} - \hat{X}_{old}\| < \epsilon$, stop.

end for

10 Iterative Solvers

Let $\tau > 0$ be step size.

Algorithm 2 Landweber Iteration

Require: $w^{(0)}$ **for** $k \leftarrow 0, 1, \dots$ **do** $w^{(k+1)} \leftarrow w^{(k)} - \tau X^T(Xw^{(k)} - y)$ **end for**

10.1 Gradient Descent/Landweber Iteration

This algorithm takes a step in the direction of negative gradient of each iteration of the objective function $f(w) = \|Xw - y\|_2^2$. Notice that

$$\begin{aligned} f(w) &= \|Xw - y\|_2^2 \\ \nabla_w f(w) &= \nabla_w (Xw - y)^T (Xw - y) \\ &= \nabla_w w^T X^T X y - 2w^T X^T y + y^T y \\ &= 2X^T Xw - 2X^T y \\ &= 2X^T (Xw - y) \end{aligned}$$

Thus the new iterate equals the old iterate plus a step in direction of negative gradient.

Claim 2. (*Convergence of Landweber Iteration*)

Proof. We want to show that

$$\|Xw^{(k+1)} - y\|_2^2 \leq \|Xw^{(k)} - y\|_2^2 \quad (49)$$

Recall that the iteration is given by $w^{(k+1)} = w^{(k)} - \tau X^T(Xw^{(k)} - y)$. Then

$$\begin{aligned} \|Xw^{(k+1)} - y\|_2^2 &= \|X(w^{(k)} - \tau X^T(Xw^{(k)} - y)) - y\|_2^2 \\ &= \|Xw^{(k)} - y - \tau XX^T(Xw^{(k)} - y)\|_2^2 \\ &= \|Xw^{(k)} - y\|_2^2 + \tau^2 \|XX^T(Xw^{(k)} - y)\|_2^2 - 2\tau(Xw^{(k)} - y)^T XX^T(Xw^{(k)} - y) \end{aligned}$$

Now observe that

$$\|XX^T(Xw^{(k)} - y)\|_2^2 \leq \|X\|_{op}^2 \|X^T(Xw^{(k)} - y)\|_2^2 \quad (50)$$

Therefore

$$\begin{aligned}\|Xw^{(k+1)} - y\|_2^2 &\leq \|Xw^{(k)} - y\|_2^2 + \tau \left(\tau \|X\|_{op}^2 \|X^T(Xw^{(k)} - y)\|_2^2 - 2 \|X^T(Xw^{(k)} - y)\|_2^2 \right) \\ &= \|Xw^{(k)} - y\|_2^2 + \tau \|X^T(Xw^{(k)} - y)\|_2^2 (\tau \|X\|_{op}^2 - 2)\end{aligned}$$

Thus, if $(\tau \|X\|_{op}^2 - 2) < 0$, then $\|Xw^{(k+1)} - y\|_2^2 \leq \|Xw^{(k)} - y\|_2^2$. Therefore, for convergence we require that

$$0 < \tau < \frac{2}{\|X\|_{op}^2} \quad (51)$$

in order to ensure convergence. Under this condition, then

$$w^{(k)} \rightarrow (X^T X)^{-1} X^T y \quad (52)$$

□

The above proof made use of the following claim.

Claim 3. (Bound on 2-norm of matrix-vector product) Let X be matrix and w a vector (conformable). Then

$$\|Xw\|_2 \leq \|X\|_{op} \|w\|_2 \quad (53)$$

Proof. Recall: **The 2-norm (here, operator norm) of a matrix is it's largest singular value.**

$$\|X\|_{op} = \|X\|_2 = \sigma_{max} \quad (54)$$

Suppose $X = U\Sigma V^T$. Then

$$\begin{aligned}\|Xw\|_2 &= \|U\Sigma V^T w\|_2 \\ &= \|\Sigma V^T w\|_2 && (U \text{ orthonormal, preserves norms}) \\ &= \left(\sum_i (\sigma_i (V^T w)_i)^2 \right)^{\frac{1}{2}} \\ &\leq \sigma_{max} \left(\sum_i (V^T w)_i^2 \right)^{\frac{1}{2}} \\ &= \sigma_{max} \|V^T w\|_2 && (\text{definition of norm}) \\ &= \sigma_{max} \|w\|_2 && (V, V^T \text{ orthonormal, preserves norms})\end{aligned}$$

□

11 Regularized Regression

11.1 Proximal Gradient Algorithm

Algorithm 3 Proximal Gradient Algorithm: $\arg \min_w \|y - Xw\|_2^2 + \lambda r(w)$

Require: Initial $w^{(0)}$

for $k \leftarrow 0, 1, \dots$ **do**

$z^{(k)} \leftarrow w^{(k)} - \tau X^T (Xw^{(k)} - y)$

▷ grad descent step

$w^{(k+1)} \leftarrow \arg \min_w \|z^{(k)} - w\|_2^2 + \lambda \tau r(w)$

▷ regularization step

end for

11.2 LASSO (Least absolute selection and shrinkage operator)

LASSO solve the following problem

$$\hat{w}_L = \arg \min_w \|w\|_1 \text{ subject to } \|y - Xw\| < \epsilon \quad (55)$$

which is equivalent to

$$\hat{w}_L = \arg \min_w \|y - Xw\|_2^2 + \lambda \|w\|_1 \quad (56)$$

In the figure below, the rhombuses and circles show the locus of points for which the weight vector has a particular norm (in the L_1 and L_2 norms respectively). More precisely, they are $\{w : \|w\|_1 = \tau_1\}$ and $\{w : \|w\|_2 = \tau_2\}$. The red line is $\{w : y = Xw\}$.

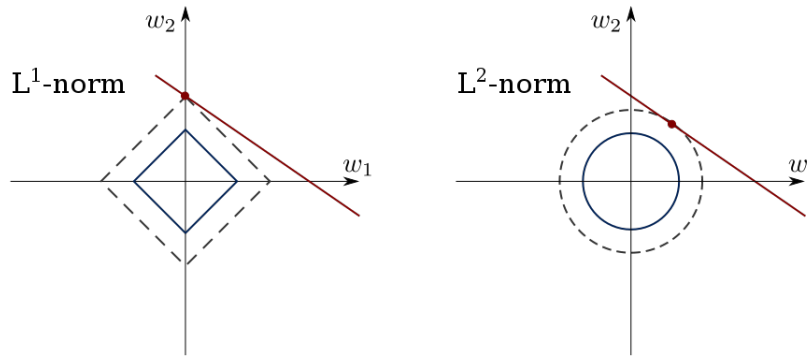


Figure 4: Weight vector with Lasso vs Ridge Regression

Example 10. ($r(w) = \|w\|_1 = \sum_i |w_i|$)

$$\begin{aligned}\hat{w} &= \arg \min_w \|z - w\|_2^2 + \lambda \tau \|w\|_1 \\ &= \arg \min_w \sum_i \left((z_i - w_i)^2 + \lambda \tau |w_i| \right)\end{aligned}\quad (\text{separable})$$

Thus we can solve a problem for each w_i .

$$\hat{w}_i = \arg \min_{w_i} (z_i - w_i)^2 + \lambda \tau |w_i| \quad (57)$$

where $\lambda, \tau > 0$.

1. Case 1: $z_i > 0$. Then $\hat{w}_i \geq 0$. Thus

$$\begin{aligned}\hat{w}_i &= \arg \min_{w_i} (z_i - w_i)^2 + \lambda \tau w_i \\ \frac{d}{dw_i}(\text{obj}) &= -2(z_i - w_i) + \lambda \tau \\ w_i &= z_i - \frac{\lambda \tau}{2}\end{aligned}$$

Thus

- (a) If $z_i > \frac{\lambda \tau}{2}$, then $\hat{w}_i = z_i - \frac{\lambda \tau}{2}$.
- (b) If $z_i < \frac{\lambda \tau}{2}$, then $\hat{w}_i = 0$.

In sum

$$\hat{w}_i = \left(z_i - \frac{\lambda \tau}{2} \right)_+ \quad (58)$$

2. Case 2: $z_i < 0$. Then $\hat{w}_i \leq 0$. Similarly,

$$w_i = z_i + \frac{\lambda \tau}{2} \quad (59)$$

Thus

- (a) If $z_i < -\frac{\lambda \tau}{2}$, then $\hat{w}_i = z_i + \frac{\lambda \tau}{2}$.
- (b) If $z_i > -\frac{\lambda \tau}{2}$, then $\hat{w}_i = 0$.

In sum

$$\hat{w}_i = - \left(|z_i| - \frac{\lambda \tau}{2} \right)_+ \quad (60)$$

We can combine these two cases to get that

$$\hat{w}_i = - \left(|z_i| - \frac{\lambda\tau}{2} \right)_+ \text{sign}(z_i) \quad (61)$$

we call the $\text{SoftThreshold}(z_i, \frac{\lambda\tau}{2})$. The figure below shows how \hat{w}_i depends on z_i .

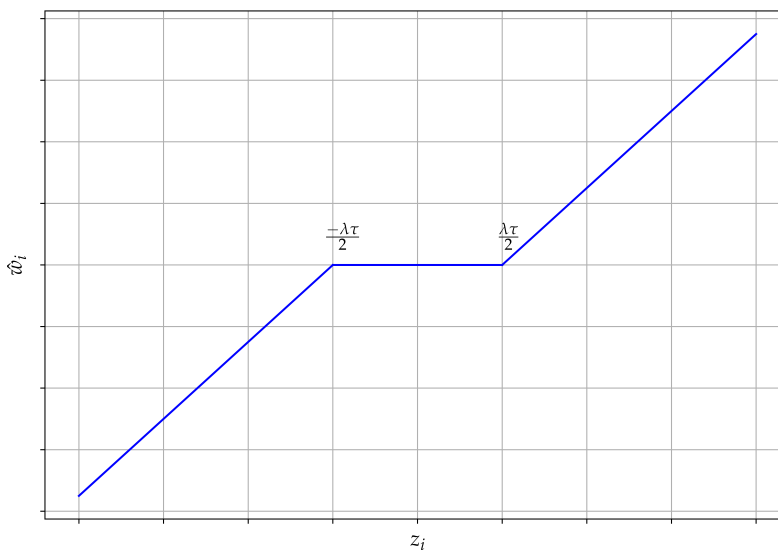


Figure 5: Soft Thresholding

12 Convexity and Support Vector Machines

12.1 Convexity

Example 11. (Convex function lies above tangent lines) Suppose $l(w)$ is a convex function. Then informally, $l(w) \geq \text{tangent}$ and w . More formally,

$$l(u) \geq l(w) + (u - w)^T \nabla l(w) \quad (62)$$

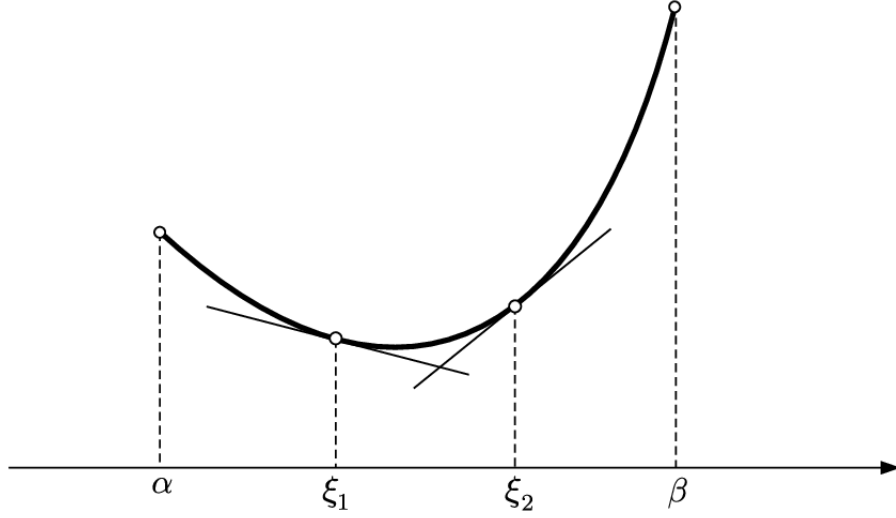


Figure 6: A convex function lies entirely above tangent lines

Definition 9. (Subgradient) A subgradient v of a convex function l is as a vector satisfying

$$l(u) \geq l(w) + (u - w)^T v \quad (63)$$

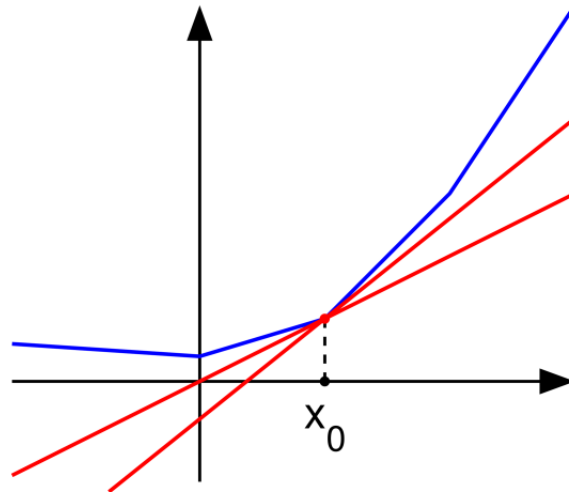


Figure 7: Subgradients

12.2 Support Vector Machines

When we use the classification rule $\hat{y} = \text{sign}(w^T x)$, our goal is thus to choose a w such that $\hat{y} = \text{sign}(w^T x)$ as often as possible. However, when we use least squares we do *not* actually minimize the number of mistakes. Minimizing the number of mistakes can be

written as minimizing the following sum of indicator variables:

$$\sum_{i=1}^n I_{y_i \neq \text{sign}(w^T x_i)} \quad (64)$$

LS actually minimizes

$$\sum_{i=1}^n (y_i - w^T x_i)^2 \quad (65)$$

We want to choose a convex function that mimics the ideal loss. We will use hinge loss, which is defined by

$$l(w) = \sum_{i=1}^n (1 - y_i x_i^T w)_+ \quad (66)$$

where

$$(a)_+ = \begin{cases} a & a > 0 \\ 0 & \text{otherwise} \end{cases} \quad (67)$$

Definition 10. (*Support Vector Machine*) If we minimize

$$\sum_{i=1}^n (1 - y_i x_i^T w)_+ + \lambda \|w\|_2^2 \quad (68)$$

this is called a support vector machine.

In the figure below, the black line is the ideal loss function (not convex). The green line is the squared loss function (LS), the blue line is hinge loss, and the red line is log loss.

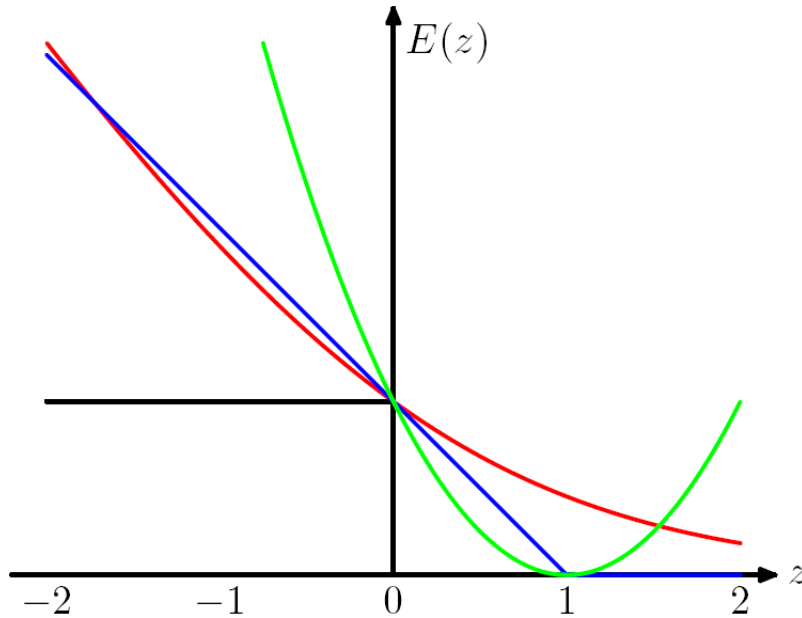


Figure 8: Loss Functions