

Linear Regression

Linear model for regression

Regression problem $\mathbb{Y} = \mathbb{R}$

ML model $a : \mathbb{X} \rightarrow \mathbb{Y}$

Family of ML models $a \in \mathcal{A}$

$$\mathcal{A} = \{a(\mathbf{x}) = w_0 + w_1x_1 + \cdots + w_Dx_D \mid w_0, w_1, \dots, w_D \in \mathbb{R}\}$$

w_0, w_1, \dots, w_D - parameters/weights of linear model

w_0 - bias

In total $D+1$ parameters. Linear models require small memory, are fast in inference and interpretable.

Linear model for regression

ML model $a(\mathbf{x}) = w_0 + \sum_{j=1}^D w_j x_j = w_0 + \langle \mathbf{w}, \mathbf{x} \rangle$

Sometimes it is convenient to augment feature vectors with constant feature =1

$$\mathbf{x} \rightarrow \tilde{\mathbf{x}} = [\mathbf{x}, 1]$$

$$a(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + w_0 = \langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}} \rangle \quad \tilde{\mathbf{w}} = [\mathbf{w}, w_0]$$

Example: apartment cost prediction

x – apartment, y – its cost

$$a(\boldsymbol{x}) = w_0 + w_1 * (\textit{area}) + w_2 * (\textit{number_of_rooms}) + \\ + w_3 * (\textit{distance_to_train_station}) + \dots$$

Important property of linear models: each feature influences independently on prediction result

Sometimes it is reasonable to consider different combinations of initial features and feature transformations

OHE for nominal features

Linear models can work only with numerical features

Let's consider nominal feature “apartment district”

$$x \in C, \quad C = \{c_1, \dots, c_m\}$$

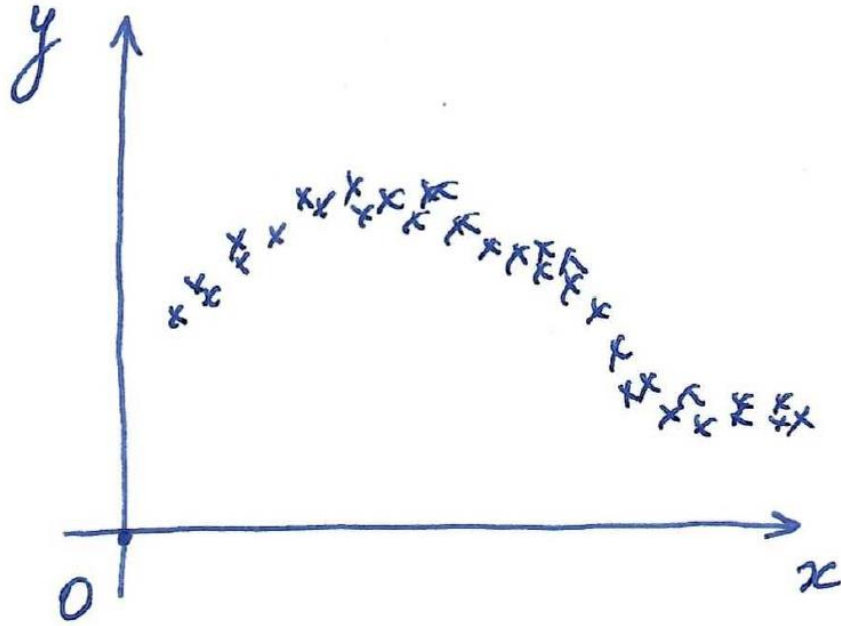
One-hot-encoding

$$x \rightarrow [b_1(x), b_2(x), \dots, b_m(x)]$$

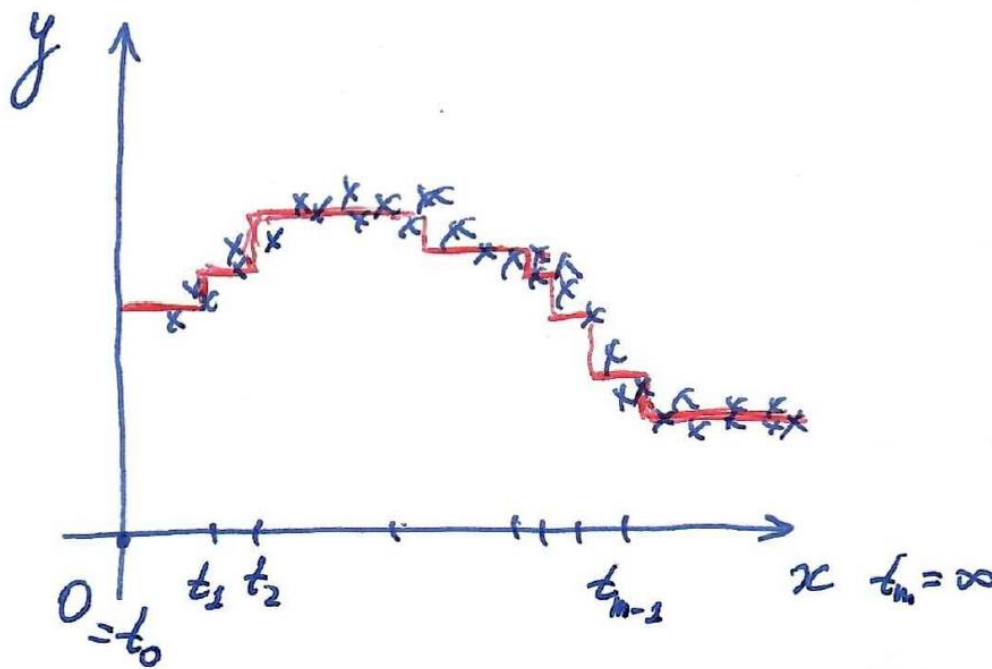
$$b_j(x) = [x = c_j] \in \{0, 1\}$$

$$a(x) = w_0 + w_1[x = c_1] + w_2[x = c_2] + \dots$$

Binarization of numerical features



Binarization of numerical features



$$x \rightarrow [b_1(x), \dots, b_m(x)]$$

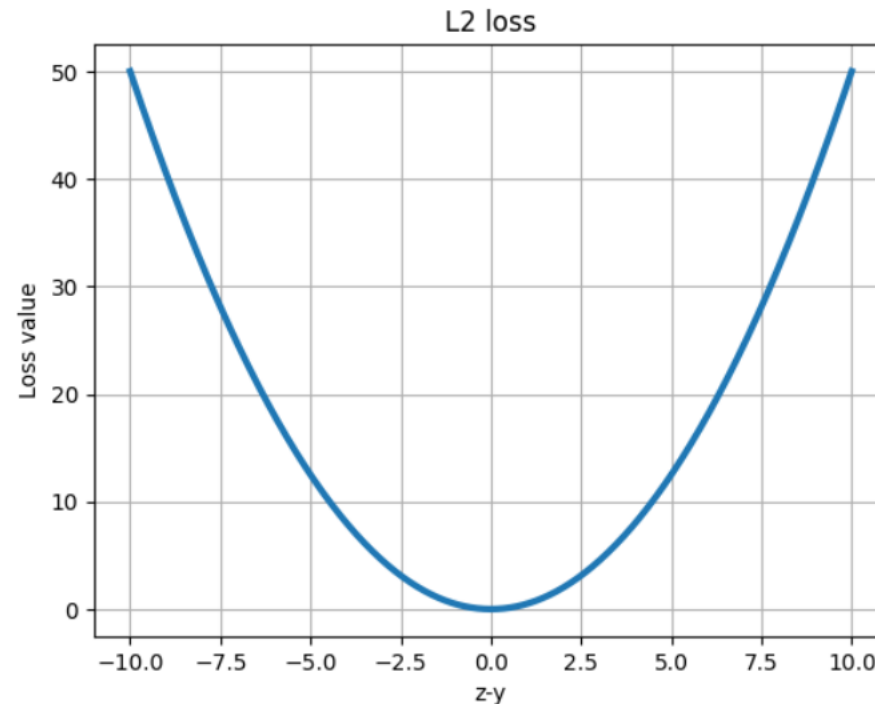
$$b_j(x) = [t_{j-1} \leq x < t_j]$$

For linear models it is important to prepare a reasonable set of features from original ones!

Loss functions for regression

Squared error loss / Quadratic loss / L2 loss $L(y, z) = (y - z)^2$

Mean squared error $MSE(a, X) = \frac{1}{N} \sum_{i=1}^N (a(\mathbf{x}_i) - y_i)^2$



Loss functions for regression

Squared error loss / Quadratic loss / L2 loss $L(y, z) = (y - z)^2$

Mean squared error $MSE(a, X) = \frac{1}{N} \sum_{i=1}^N (a(\mathbf{x}_i) - y_i)^2$

Root mean squared error $RMSE(a, X) = \sqrt{MSE(a, X)}$

Coefficient of determination

$$R^2(a, X) = 1 - \frac{\sum_{i=1}^N (a(\mathbf{x}_i) - y_i)^2}{\sum_{i=1}^N (y_i - \bar{y})^2}$$

For ideal $a(x)$ $R^2=1$, for constant $a(x)$ $R^2=0$

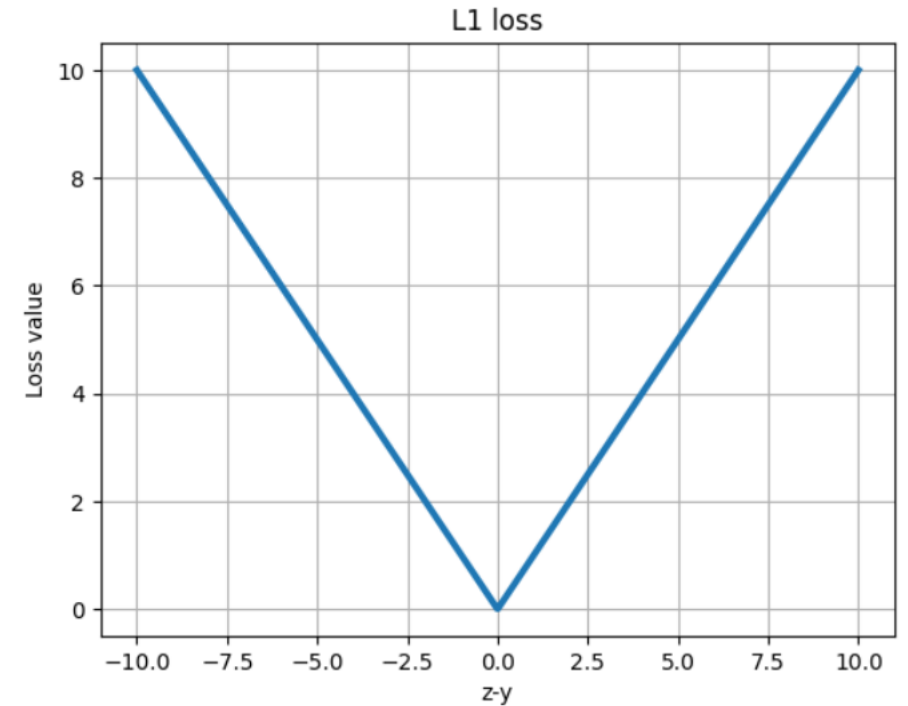
Mean Absolute Error loss

L1 loss $L(y, z) = |y - z|$

$$MAE(a, X) = \frac{1}{N} \sum_{i=1}^N |a(\mathbf{x}_i) - y_i|$$

MAE is non-differentiable

MAE is more robust to outliers

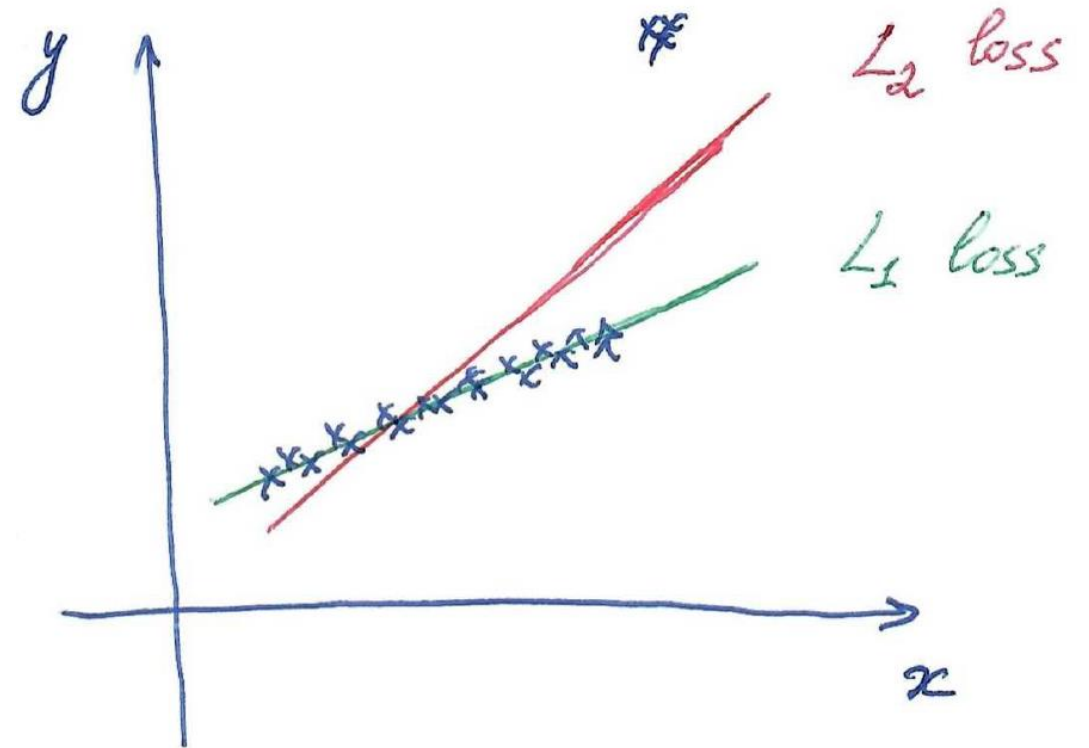


Mean Absolute Error loss

L1 loss $L(y, z) = |y - z|$

$$MAE(a, X) = \frac{1}{N} \sum_{i=1}^N |a(\mathbf{x}_i) - y_i|$$

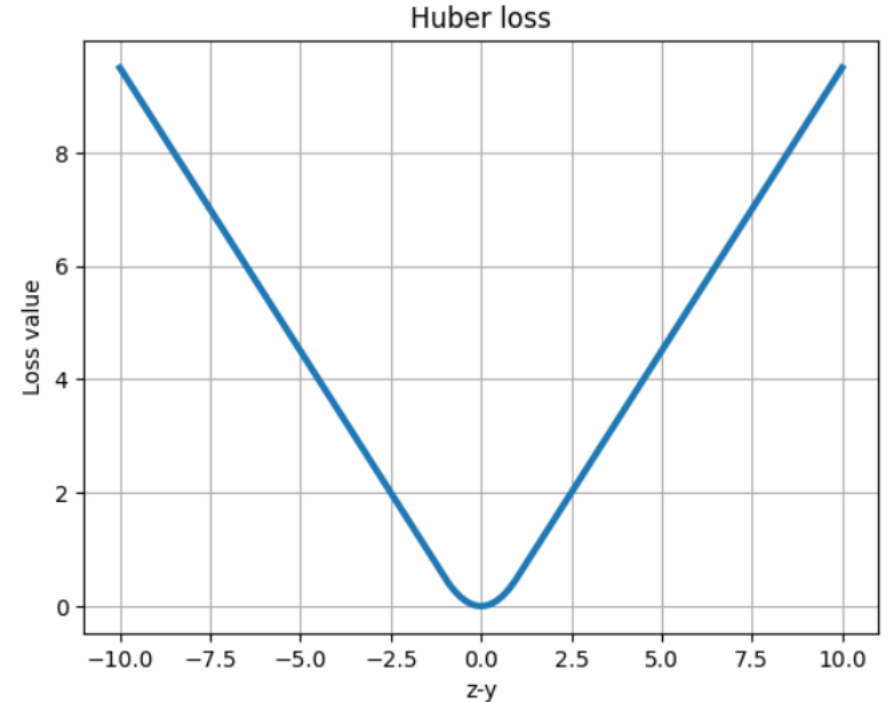
y	z	$(y - z)^2$	$ y - z $
1	2	1	1
1000	2	996004	998
1	3	4	2
1000	3	994009	997



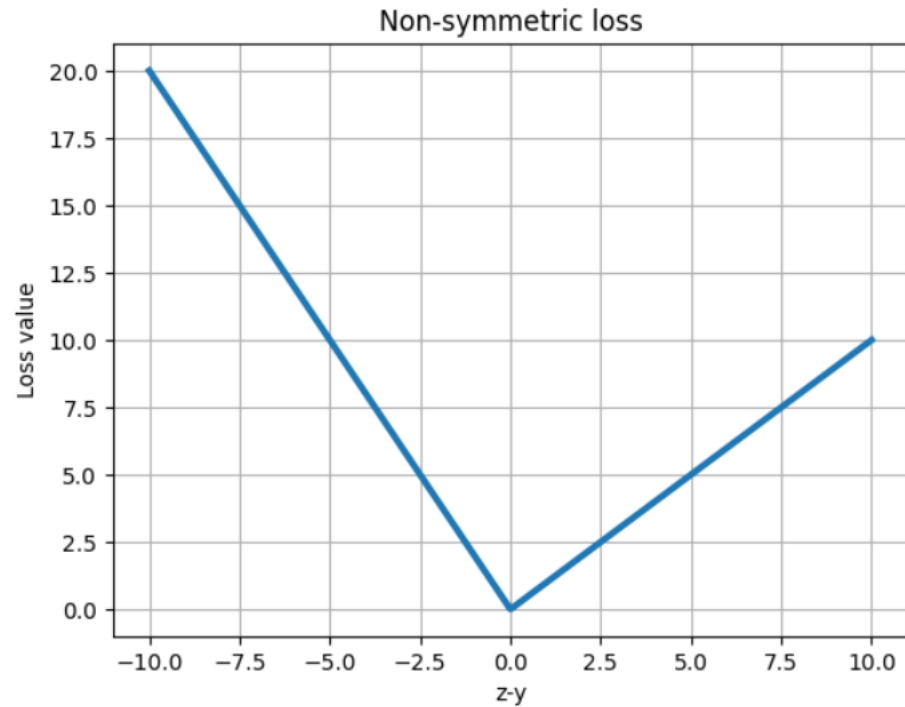
Huber loss

$$L(y, z) = \begin{cases} \frac{1}{2}(y - z)^2, & \text{if } |y - z| \leq \delta, \\ \delta|y - z| - \frac{1}{2}\delta^2, & \text{otherwise.} \end{cases}$$

Huber loss is continuously differentiable everywhere



Non-symmetric loss



Loss function should be chosen according to business requirements

Other losses

Mean Squared Logarithmic Error (MSLE)

$$L(y, z) = (\log(z + 1) - \log(y + 1))^2$$

Mean Absolute Percentage Error (MAPE)

$$L(y, z) = \left| \frac{y - z}{y} \right|$$

And many others...

Linear regression training

$$Q(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N (\langle \mathbf{w}, \mathbf{x}_i \rangle - y_i)^2 \rightarrow \min_{\mathbf{w}}$$

$$X = \begin{bmatrix} x_{1,1} & \dots & x_{1,D} \\ x_{2,1} & \dots & x_{2,D} \\ \vdots & & \vdots \\ x_{N,1} & \dots & x_{N,D} \end{bmatrix} \in \mathbb{R}^{N \times D} \quad \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix} \in \mathbb{R}^D \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \in \mathbb{R}^N$$

$$z_i = \langle \mathbf{w}, \mathbf{x}_i \rangle - y_i$$

$$\mathbf{z} = X\mathbf{w} - \mathbf{y}$$

$$Q(w) = \frac{1}{N} \sum_{i=1}^N z_i^2 = \frac{1}{N} \|\mathbf{z}\|_2^2 = \frac{1}{N} \|X\mathbf{w} - \mathbf{y}\|_2^2$$

Linear regression training

$$Q(\mathbf{w}) = \frac{1}{N} \|X\mathbf{w} - \mathbf{y}\|_2^2 \rightarrow \min_{\mathbf{w}}$$

$$\nabla_{\mathbf{w}} Q(\mathbf{w}) = \mathbf{0} \Rightarrow \mathbf{w}_{opt} = (X^T X)^{-1} X^T \mathbf{y}$$

\mathbf{w}_{opt} is solution of linear system $X^T X \mathbf{w}_{opt} = X^T \mathbf{y}$

Properties of the matrix $X^T X$:

- 1) Quadratic and symmetric: $(X^T X)^T = X^T X$
- 2) Non-negatively defined: $\mathbf{u}^T X^T \underbrace{X\mathbf{u}}_{\mathbf{v}} = \mathbf{v}^T \mathbf{v} \geq 0$

Recall: matrix A is non-negatively definite, if $\mathbf{u}^T A \mathbf{u} \geq 0 \forall \mathbf{u}$

Equivalently: all eigenvalues are non-negative.

Linear regression training

\mathbf{w}_{opt} is solution of linear system $X^T X \mathbf{w}_{opt} = X^T \mathbf{y}$

For analytical solution the matrix $X^T X$ should be non-singular.

$$\text{rank}(X^T X) \leq \min(\text{rank}(X^T), \text{rank}(X))$$

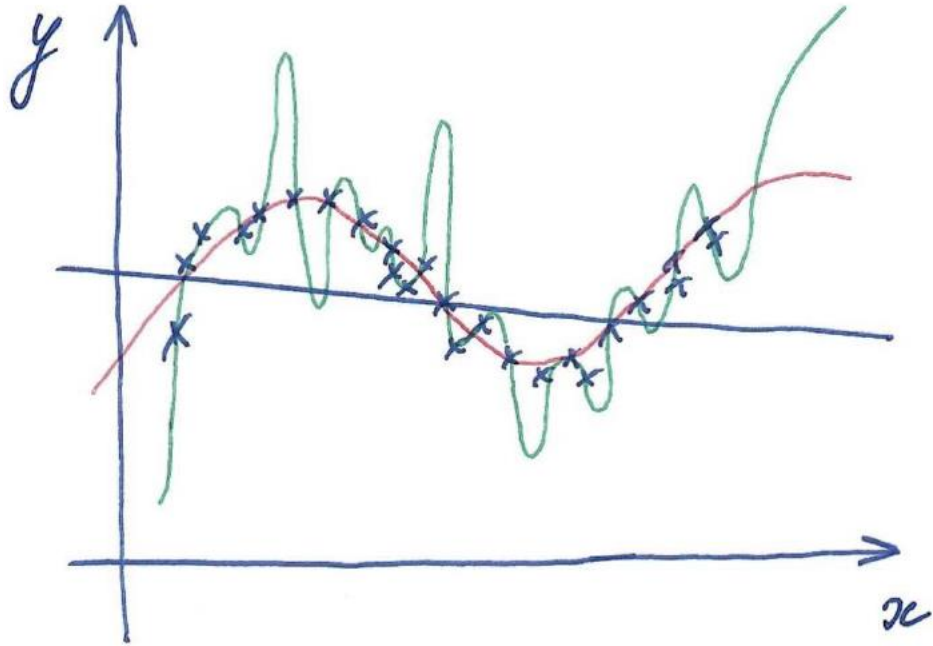
So the matrix is singular if $N < D$ or there are duplicating features

Computational complexity of analytical solution:

$O(ND^2)$ for computing $X^T X$

$O(D^3)$ for linear system solving

Underfitting and overfitting



$$a(x) = w_0 + w_1x$$

$$a(x) = w_0 + w_1x + w_2x^2 + w_3x^3$$

$$a(x) = w_0 + w_1x + \dots + w_nx^n$$

Blue model is too simple and underfitted, green model is too complicated and overfitted, red model has good generalization ability

Regularization

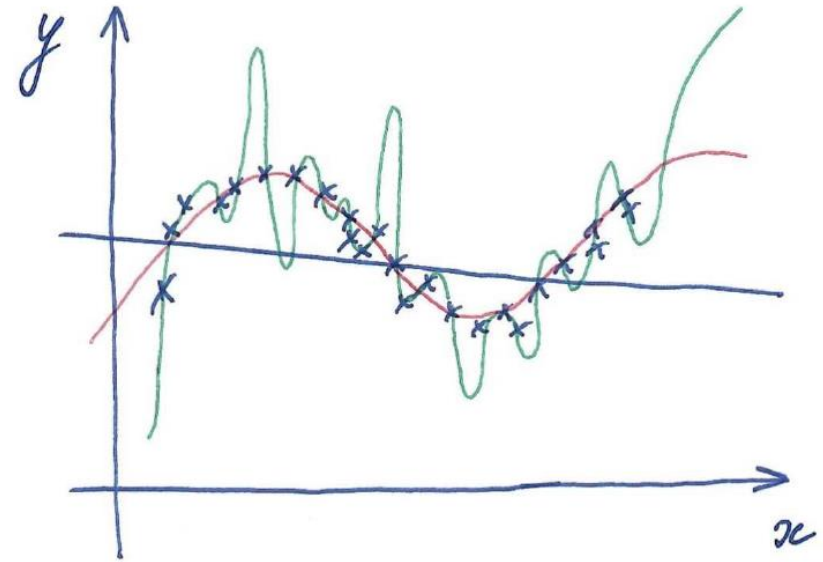
Overfitted model usually has oscillating effect – small change of inputs may lead to large change of outputs

$$a(\mathbf{x}, \mathbf{w}) = w_0 + w_1x_1 + \cdots + w_Dx_D$$

For linear models oscillating happens with high absolute values of some weights

Idea: let's additionally require during training that our ML model should have weights with low absolute values

$$Q(\mathbf{w}, X, \mathbf{y}) + \lambda \|\mathbf{w}\|_2^2 \rightarrow \min_{\mathbf{w}}, \lambda > 0$$



Regularization

$$F(\mathbf{w}) = Q(\mathbf{w}, X, \mathbf{y}) + \lambda \|\mathbf{w}\|_2^2 = \frac{1}{N} \|X\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 \rightarrow \min_{\mathbf{w}}$$

This is known as L2 regularization, weight decay, ridge regularization

This minimization problem can be solved analytically

$$\mathbf{w}_{opt} = (X^T X + \lambda N I)^{-1} X^T \mathbf{y}$$

Regularization

Non-regularized solution $\mathbf{w}_{opt} = (X^T X)^{-1} X^T \mathbf{y}$

This solution can't be computed if $N < D$ or there are duplicating features

Regularized solution $\mathbf{w}_{opt} = (X^T X + \lambda N I)^{-1} X^T \mathbf{y}$

This solution can always be computed because the matrix $X^T X + \lambda N I$ is strictly positively definite

Important: bias parameter is usually not regularized!

Eigenvalue Decomposition

Let's consider arbitrary symmetric matrix $A^T = A \in \mathbb{R}^{n \times n}$
and all its eigenvalues $\lambda_1, \dots, \lambda_n$ and eigenvectors $\mathbf{q}_1, \dots, \mathbf{q}_n$

Eigenvalue decomposition in matrix form: $A = Q\Lambda Q^T$

$$\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n) \quad Q = [\mathbf{q}_1, \dots, \mathbf{q}_n], \quad QQ^T = I$$

$$X^T X = Q\Lambda Q^T \quad X^T X + \lambda N I = Q\Lambda Q^T + \lambda N Q Q^T = Q[\Lambda + \lambda N I]Q^T$$

$$\text{eigenvalue}(X^T X + \lambda N I) = \text{eigenvalue}(X^T X) + \lambda N$$

Choosing regularization coefficient

$$F(\mathbf{w}, \lambda) = \frac{1}{N} \|X\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 \rightarrow \min_{\mathbf{w}, \lambda} \Rightarrow \lambda_{opt} = 0$$

Regularization coefficient can't be optimized on the training set. Parameters with this property are called hyperparameters.

Hyperparameter can be chosen by brute force – estimating cross-validation (CV) value of ML model for some set of values

(1/N) coefficient in loss function helps making optimal value of regularization parameter less sensitive to the training set size N.

Underfitting and overfitting

In case of overfitting increase regularization coefficient

In case of underfitting decrease regularization coefficient or increase ML model complexity

L1 regularization / Lasso regression

$$F(\mathbf{w}) = \frac{1}{N} \|X\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_1 = \frac{1}{N} \|X\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \sum_{j=1}^D |w_j| \rightarrow \min_{\mathbf{w}}$$

L1 norm penalizes high absolute values of weights and thus prevents overfitting.

L1 norm allows to find sparse solution when some of weights w_j exactly equal zero. In such a way we may find relevant/irrelevant features.

$$a(\mathbf{x}) = w_0 + w_1 * (area) + w_2 * (number_of_rooms) + \\ + w_3 * (distance_to_train_station) + \dots$$

L1 regularization / Lasso regression

From optimization theory it is known that the optimization problem

$$f(\boldsymbol{w}) + \lambda g(\boldsymbol{w}) \rightarrow \min_{\boldsymbol{w}}$$

is equivalent to the following constrained optimization problem:

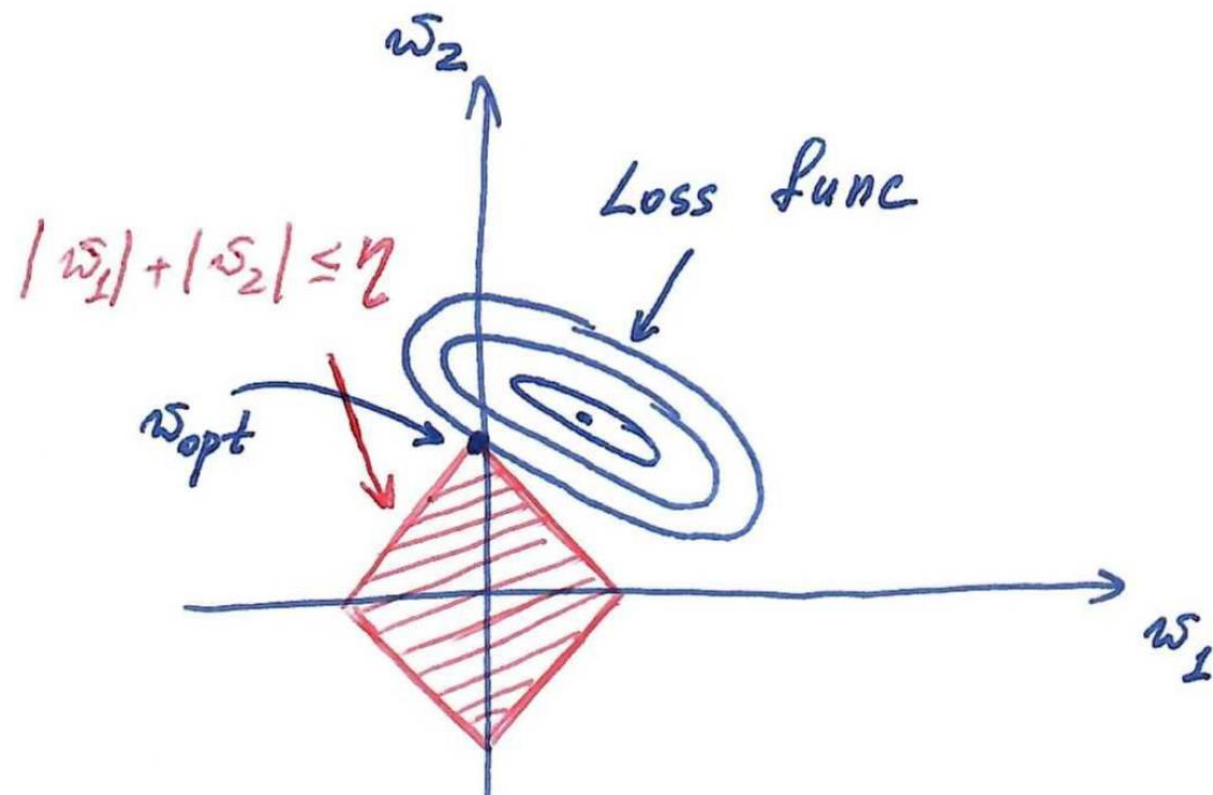
$$\begin{aligned} f(\boldsymbol{w}) &\rightarrow \min_{\boldsymbol{w}}, \\ \text{s.t. } g(\boldsymbol{w}) &\leq \eta \end{aligned}$$

L1 regularization / Lasso regression

$$f(\mathbf{w}) + \lambda \|\mathbf{w}\|_1 \rightarrow \min_{\mathbf{w}}$$

$$f(\mathbf{w}) \rightarrow \min_{\mathbf{w}},$$

$$\text{s.t. } \|\mathbf{w}\|_1 \leq \eta$$

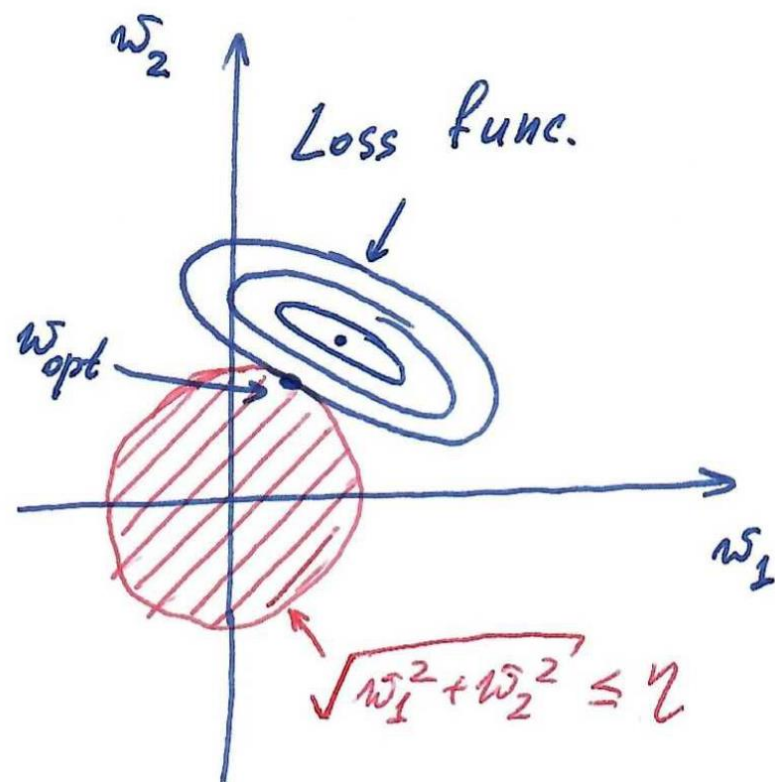


L2 regularization / Ridge regression

$$f(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2 \rightarrow \min_{\mathbf{w}}$$

$$f(\mathbf{w}) \rightarrow \min_{\mathbf{w}},$$

$$\text{s.t. } \|\mathbf{w}\|_2^2 \leq \eta$$



Elastic Net regularization

$$\frac{1}{N} \|X\boldsymbol{w} - \boldsymbol{y}\|_2^2 + \lambda_1 \|\boldsymbol{w}\|_2^2 + \lambda_2 \|\boldsymbol{w}\|_1 = \frac{1}{N} \|X\boldsymbol{w} - \boldsymbol{y}\|_2^2 + \lambda_1 ((1 - \lambda_2) \|\boldsymbol{w}\|_2^2 + \lambda_2 \|\boldsymbol{w}\|_1)$$

Matrix/vector differentiation

Vector expressions

In mathematical expressions a column-wise vector notation is convenient

Scalar

$$\alpha \in \mathbb{R} = \mathbb{R}^1 = \mathbb{R}^{1 \times 1}$$



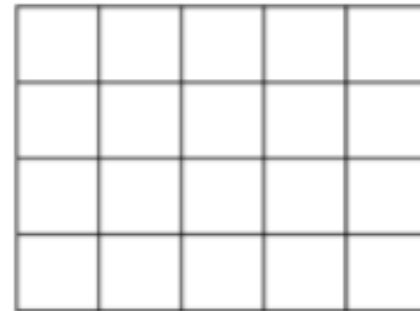
Vector

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{bmatrix} \in \mathbb{R}^n = \mathbb{R}^{n \times 1}$$



Matrix

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1m} \\ a_{21} & a_{22} & \dots & a_{2m} \\ \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nm} \end{bmatrix} \in \mathbb{R}^{n \times m}$$



Vector expressions

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^n$$

$$\mathbf{x}, \mathbf{y} \in \mathbb{R}^n \quad \langle \mathbf{x}, \mathbf{y} \rangle = \sum_{i=1}^n x_i y_i = \mathbf{x}^T \mathbf{y}$$

$$\mathbf{x} \mathbf{y}^T \in \mathbb{R}^{n \times n} \quad (\mathbf{x} \mathbf{y}^T)_{ij} = x_i y_j$$

$$\|\mathbf{x}\|_2^2 = \sum_{i=1}^n x_i^2 = \mathbf{x}^T \mathbf{x}$$

Matrix expressions

$$\begin{aligned} A, B \in \mathbb{R}^{n \times m} \quad \langle A, B \rangle &= \sum_{i,j=1}^{n,m} A_{ij} B_{ij} = \text{tr}(A^T B) = \\ &= \sum_{j=1}^m (A^T B)_{jj} = \sum_{j=1}^m \sum_{i=1}^n (A^T)_{ji} B_{ij} = \sum_{j=1}^m \sum_{i=1}^n A_{ij} B_{ij} \end{aligned}$$

$$\|A\|_F^2 = \sum_{i,j} A_{ij}^2 = \text{tr}(A^T A)$$

Properties of transposition: $(A + B)^T = A^T + B^T$

$$(AB)^T = B^T A^T$$

Summation and product

Summation of two vector/matrix objects $x + y$
is valid only if two objects have the same size

Product of two objects AB
is valid only if the last dimension of the first object coincides with
the first dimension of the second object.

Exception: multiplication by scalar is allowed. αA

In general, the order of arguments in product is important because

$$AB \neq BA$$

Summation and product

$$\boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^n, \boldsymbol{z} \in \mathbb{R}^m, \boldsymbol{A} \in \mathbb{R}^{n \times m}, b \in \mathbb{R}$$

Correct examples:

1. $\boldsymbol{x} + \boldsymbol{y}$
2. $\boldsymbol{x}^T \boldsymbol{y} + b$
3. $\boldsymbol{A} \boldsymbol{z}$
4. $\boldsymbol{A} + \boldsymbol{x} \boldsymbol{z}^T$
5. $\boldsymbol{A} \boldsymbol{A}^T$
6. $b \boldsymbol{A}$

Incorrect examples:

1. $\boldsymbol{x} + \boldsymbol{z}$
2. $\boldsymbol{x} + b$
3. $\boldsymbol{x} + \boldsymbol{y}^T$
4. $\boldsymbol{A} \boldsymbol{x}$
5. $\boldsymbol{A} \boldsymbol{A}$
6. $\boldsymbol{A} + \boldsymbol{x} \boldsymbol{y}^T$

Inversion (division)

It is possible to multiply both parts of equation either by inversed scalar α^{-1} or by inversed matrix A^{-1} , where A is quadratic and non-degenerate matrix.

Division by vector x^{-1} is not possible!

$$x, y \in \mathbb{R}^n, A \in \mathbb{R}^{n \times n}, \alpha \in \mathbb{R}$$

Correct examples:

$$1. \alpha x = y \Rightarrow x = \alpha^{-1}y$$

$$2. Ax = y \Rightarrow x = A^{-1}y$$

Incorrect examples:

$$1. x^T y = \alpha \Rightarrow y = x^{-T} \alpha$$

$$2. X^T X w = X^T y \Rightarrow X w = X^{-T} X^T y = y$$

Manipulating with vector expression

$$\begin{aligned} Q(\mathbf{w}) &= \frac{1}{N} \|X\mathbf{w} - \mathbf{y}\|_2^2 = \frac{1}{N} (X\mathbf{w} - \mathbf{y})^T (X\mathbf{w} - \mathbf{y}) = \frac{1}{N} ((X\mathbf{w})^T - \mathbf{y}^T) (X\mathbf{w} - \mathbf{y}) = \\ &= \frac{1}{N} (\mathbf{w}^T X^T - \mathbf{y}^T) (X\mathbf{w} - \mathbf{y}) = \frac{1}{N} (\mathbf{w}^T X^T X \mathbf{w} - \mathbf{y}^T X \mathbf{w} - \mathbf{w}^T X^T \mathbf{y} + \mathbf{y}^T \mathbf{y}) \end{aligned}$$

Check yourself by dimensionality! $Q(\mathbf{w}) \in \mathbb{R}$ $\underbrace{\mathbf{y}^T}_{1 \times N} \underbrace{X}_{N \times D} \underbrace{\mathbf{w}}_{D \times 1} \in \mathbb{R}$

$$\underbrace{\mathbf{y}^T X \mathbf{w}}_{\in \mathbb{R}} = (\mathbf{y}^T X \mathbf{w})^T = \mathbf{w}^T X^T \mathbf{y}$$

$$Q(\mathbf{w}) = \frac{1}{N} (\mathbf{w}^T X^T X \mathbf{w} - 2\mathbf{w}^T X^T \mathbf{y} + \mathbf{y}^T \mathbf{y})$$

Manipulating with vector expression

$$\mathbf{u}, \mathbf{v} \in \mathbb{R}^n \quad \|\mathbf{u}\mathbf{v}^T\|_F^2 = \text{tr}((\mathbf{u}\mathbf{v}^T)^T \mathbf{u}\mathbf{v}^T) = \text{tr}(\mathbf{v}\mathbf{u}^T \mathbf{u}\mathbf{v}^T)$$

Circular property of
trace operation:

$$\text{tr}(ABC) = \text{tr}(CAB) = \text{tr}(BCA)$$

$$\text{tr}(\mathbf{v}\mathbf{u}^T \mathbf{u}\mathbf{v}^T) = \text{tr}(\underbrace{\mathbf{v}^T \mathbf{v}}_{\in \mathbb{R}} \underbrace{\mathbf{u}^T \mathbf{u}}_{\in \mathbb{R}}) = \mathbf{v}^T \mathbf{v} \mathbf{u}^T \mathbf{u} = \|\mathbf{v}\|_2^2 \|\mathbf{u}\|_2^2$$

Differential calculus in scalar case

$$f : \mathbb{R} \rightarrow \mathbb{R} \qquad f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

Table of standard
derivatives:

$$(x^n)' = nx^{n-1}, \quad (\sin(x))' = \cos(x), \quad (\exp(x))' = \exp(x), \dots$$

Differentiation rules: $(f+g)' = f' + g'$, $(fg)' = f'g + fg'$, $(f(g))' = f'(g')$, \dots

Differential calculus in vector/matrix case

$$f(X) = \text{tr}(AX^{-1}B) : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}$$

$$\frac{\partial f(X)}{\partial X} = \left\{ \frac{\partial f}{\partial X_{ij}} \right\}_{i,j=1}^{n,n} = \frac{\partial}{\partial X} \text{tr} \left(\underbrace{AX^{-1}B}_{\text{2D matrix}} \right) = \text{tr} \left[\underbrace{\frac{\partial}{\partial X} AX^{-1}B}_{\text{4D tensor}} \right]$$

Goal: Make differentiation without dealing with higher dimensional tensors in intermediate computations

Differential calculus in vector/matrix case

$f : \mathcal{U} \rightarrow \mathcal{V}$ \mathcal{U}, \mathcal{V} - set of scalars, vectors, matrices

Def. Function f is called differentiable at point x if

$f(x + h) - f(x) = df(x)[h] + \bar{o}(\|h\|)$, $df(x)[h]$ - linear function w.r.t. h

$df(x)[h] = df(x)[dx] = df$ - differential

For fixed x $df(x)[h] : \mathcal{U} \rightarrow \mathcal{V}$

For scalar functions $f : \mathbb{R} \rightarrow \mathbb{R}$ differential is just a product of derivative and increment $df = df(x)[dx] = f'(x)dx$

For functions $f : \mathbb{R}^n \rightarrow \mathbb{R}$ differential is $df = df(x)[dx] = \nabla f(x)^T dx$

Differential calculus in vector/matrix case

Differentials of
standard functions:

$$d(\mathbf{c}^T \mathbf{x}) = \mathbf{c}^T d\mathbf{x}$$

$$\begin{aligned} d(\mathbf{x}^T A \mathbf{x}) &= \mathbf{x}^T (A + A^T) d\mathbf{x} = \\ &= \{A^T = A\} = 2\mathbf{x}^T A d\mathbf{x} \end{aligned}$$

$$d\operatorname{tr}(X) = \operatorname{tr}(dX)$$

$$d\det(X) = \det(X) \operatorname{tr}(X^{-1} dX)$$

$$dX^{-1} = -X^{-1} dX X^{-1}$$

Differentiation rules:

$$dA = 0$$

$$d(AXB) = A \cdot dX \cdot B$$

$$d(X + Y) = dX + dY$$

$$d(XY) = dX \cdot Y + X \cdot dY$$

$$d(X/\varphi) = \frac{dX\varphi - Xd\varphi}{\varphi^2}, \quad \varphi \in \mathbb{R}$$

$$f(x) = h(g(x))$$

$$\begin{aligned} df(x)[dx] &= dh(g(x))[dg] = \\ &= dh(g(x))[dg(x)[dx]] \end{aligned}$$

Differential calculus in vector/matrix case

$$f(\mathbf{x}) = \frac{1}{3} \|\mathbf{x}\|_2^3 = \frac{1}{3} (\mathbf{x}^T \mathbf{x})^{\frac{3}{2}} : \mathbb{R}^n \rightarrow \mathbb{R} \quad \nabla_{\mathbf{x}} f(\mathbf{x}) - ?$$

$$f(\mathbf{x}) = h(g(\mathbf{x})), \quad g(\mathbf{x}) = \mathbf{x}^T \mathbf{x}, \quad h(g) = \frac{1}{3} g^{\frac{3}{2}}$$

$$dg = dg(\mathbf{x})[d\mathbf{x}] = 2\mathbf{x}^T d\mathbf{x}$$

$$dh = dh(g)[dg] = h'(g)dg = \frac{1}{3} \cdot \frac{3}{2} \cdot g^{\frac{1}{2}} dg = \frac{1}{2} g^{\frac{1}{2}} dg$$

$$df = df(\mathbf{x})[d\mathbf{x}] = dh(g(\mathbf{x}))[dg] = \frac{1}{2} g^{\frac{1}{2}} 2\mathbf{x}^T d\mathbf{x} = (\mathbf{x}^T \mathbf{x})^{\frac{1}{2}} \mathbf{x}^T d\mathbf{x}$$

Differential calculus in vector/matrix case

Connection between differential and gradient:

$$f : \mathcal{U} \rightarrow \mathbb{R}$$

\mathcal{U}	\mathbb{R}	\mathbb{R}^n	$\mathbb{R}^{n \times m}$
Canonical form of df	$f'(x)dx$	$\nabla f(\mathbf{x})^T d\mathbf{x}$	$\text{tr}(\nabla f(X)^T dX)$

$$df = (\mathbf{x}^T \mathbf{x})^{\frac{1}{2}} \mathbf{x}^T d\mathbf{x} = \nabla f(\mathbf{x})^T d\mathbf{x}$$

$$\nabla f(\mathbf{x}) = (\mathbf{x}^T \mathbf{x})^{\frac{1}{2}} \mathbf{x} = \|\mathbf{x}\|_2 \mathbf{x}$$

Differential calculus in vector/matrix case

General scheme:

1. Compute differential of your function using table of standard differentials and differentiation rules
2. Transform expression for differential into canonical form
3. Extract gradient value

Differential calculus for linear regression

$$Q(\mathbf{w}) = \frac{1}{N} \|X\mathbf{w} - \mathbf{y}\|_2^2 = \frac{1}{N} (\mathbf{w}^T X^T X \mathbf{w} - 2\mathbf{w}^T X^T \mathbf{y} + \mathbf{y}^T \mathbf{y})$$

$$dQ = \frac{1}{N} (2\mathbf{w}^T X^T X d\mathbf{w} - 2d\mathbf{w}^T X^T \mathbf{y}) = \nabla Q(\mathbf{w})^T d\mathbf{w}$$

$$d\mathbf{w}^T X^T \mathbf{y} = \mathbf{y}^T X d\mathbf{w}$$

$$\nabla Q(\mathbf{w}) = \frac{1}{N} (2X^T X \mathbf{w} - 2X^T \mathbf{y}) = \frac{2}{N} X^T (X\mathbf{w} - \mathbf{y}) = \mathbf{0}$$

$$\Rightarrow X^T X \mathbf{w} = X^T \mathbf{y} \Rightarrow \mathbf{w}_{opt} = (X^T X)^{-1} X^T \mathbf{y}$$

Differential calculus for linear regression

$$Q(\boldsymbol{w}) = \frac{1}{N} \|\underbrace{X\boldsymbol{w} - \boldsymbol{y}}_{\boldsymbol{z}}\|_2^2 = \frac{1}{N} \|\boldsymbol{z}\|_2^2 = \boldsymbol{z}^T \boldsymbol{z}$$

$$dQ = \frac{1}{N} 2\boldsymbol{z}^T d\boldsymbol{z} = \frac{2}{N} \boldsymbol{z}^T (X d\boldsymbol{w}) = \frac{2}{N} (X\boldsymbol{w} - \boldsymbol{y})^T X d\boldsymbol{w}$$

$$\Rightarrow \nabla Q(\boldsymbol{w}) = \frac{2}{N} X^T (X\boldsymbol{w} - \boldsymbol{y})$$

Differential calculus for linear regression

$$Q(\mathbf{z}) = \frac{1}{N} \mathbf{z}^T \mathbf{z}, \quad \mathbf{z} = X\mathbf{w} - \mathbf{y} \qquad \nabla_X Q(X) - ?$$

$$\begin{aligned} dQ &= \frac{1}{N} 2\mathbf{z}^T d\mathbf{z} = \frac{2}{N} \mathbf{z}^T (dX\mathbf{w}) = \text{tr}(\nabla_X Q(X)^T dX) = \\ &= \frac{2}{N} \text{tr}(\mathbf{z}^T dX\mathbf{w}) = \frac{2}{N} \text{tr}(\mathbf{w}\mathbf{z}^T dX) \end{aligned}$$

$$\Rightarrow \nabla_X Q(X) = \frac{2}{N} \mathbf{z}\mathbf{w}^T = \frac{2}{N} (X\mathbf{w} - \mathbf{y})\mathbf{w}^T$$

$$X \in \mathbb{R}^{N \times D}, \quad \nabla_X Q(X) = \left\{ \frac{\partial Q(X)}{\partial X_{ij}} \right\}_{i,j=1}^{N,D} \in \mathbb{R}^{N \times D}, \quad \mathbf{z} \in \mathbb{R}^N, \quad \mathbf{w} \in \mathbb{R}^D, \quad \mathbf{z}\mathbf{w}^T \in \mathbb{R}^{N \times D}$$

Matrix differential example

$$f(X) = \det(AX^{-1}B) : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}$$

$$\begin{aligned} df &= \det(AX^{-1}B) \operatorname{tr}((AX^{-1}B)^{-1} d(AX^{-1}B)) = \\ &= \det(AX^{-1}B) \operatorname{tr}(-(AX^{-1}B)^{-1} AX^{-1} dX X^{-1} B) = \\ &= -\det(AX^{-1}B) \operatorname{tr}(X^{-1} B (AX^{-1}B)^{-1} AX^{-1} dX) \end{aligned}$$

$$\nabla_X f(X) = -\det(AX^{-1}B) X^{-T} A^T (AX^{-1}B)^{-T} B^T X^{-T}$$