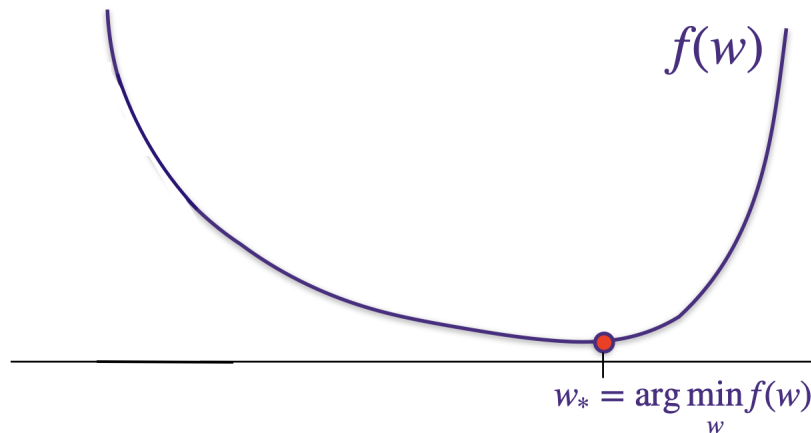


Section 04: Gradient Descent, Generalized Least Squares Regression, Convexity, MAP as Regularization

1. Gradient Descent

Like we've seen in lecture, gradient descent is an important algorithm commonly used to train machine learning models, particularly useful for when there is no closed form solution for the minimum of a loss function. Here, we'll go through short introduction to the algorithm.

Consider some function $f(w)$, which has some w_* for which $w_* = \arg \min_w f(w)$:



Let w_0 be some initial guess for the minimum of $f(w)$. Gradient descent will allow us to improve this solution.

(a) For some w that is very close to w_0 , give the Taylor series approximation for $f(w)$ starting at $f(w_0)$.

(b) Now, let us choose some $\eta > 0$ that is *very small*. With this very small η , let's assume that $w_1 = w_0 - \eta \left(\frac{df(w)}{dw} \Big|_{w=w_0} \right)$. Using your approximation from part (a), give an expression for $f(w_1)$.

(c) Given your expression for $f(w_1)$ from part (b), explain why, if η is small enough and if the function approximation is a good enough approximation, we are guaranteed to move in the “right” direction closer to the minimum w_* .

(d) Building from your answer in part (c), write a general form for the gradient descent algorithm.

2. Generalized Least Squares Regression

In class, we've seen linear regression and ridge regression. Here, we consider a problem that generalizes both of these. As a reminder, in linear regression, we seek a model that captures a linear relationship between input data and output data. The general case we consider imposes additional structure on the model.

Consider an experiment in which you have n data points $x_i \in \mathbb{R}^d$ and corresponding n observations y_i . We wish to come up with a model $\omega \in \mathbb{R}^d$ that satisfies the following properties: first, the error $\sum_{i=1}^n (x_i^\top \omega - y_i)^2$ should be small; second, we don't want small changes in training data resulting in large changes in solution; third, we want to put different weights in controlling the magnitude of different coordinates of ω . We therefore define

$$\hat{\omega}_{\text{general}} = \arg \min_{\omega} \sum_{i=1}^n (y_i - x_i^\top \omega)^2 + \lambda \sum_{i=1}^d D_{ii} \omega_i^2.$$

Here, D is a diagonal matrix, with positive entries on the diagonal. Observe that when D is the identity matrix, we recover ridge regression, and when $\lambda = 0$, we recover least squares regression. Different weights on D_{ii} cause the magnitudes of ω_i to be controlled differently.

2.1. Closed form in the general case

Deduce the closed form solution for $\hat{\omega}_{\text{general}}$. You should be comfortable with proofs in the "coordinate" form as well as the "matrix" form.

2.2. Special cases: linear regression and ridge regression

- (a) In the simple least squares case ($\lambda = 0$ above), what happens to the resulting \hat{w} if we double all the values of y_i ?
- (b) In the simple least squares case ($\lambda = 0$ above), what happens to the resulting \hat{w} if we double the data matrix $X \in \mathbb{R}^{n \times d}$?
- (c) Suppose $D = I$ (that is, it is the identity matrix). That is, this is the *ridge* regression setting. Explain why $\lambda > 0$ ensures that the solution exists and the matrix can be inverted.

3. Convexity

Convexity is defined for both sets and functions. For today we'll focus on discussing the convexity of functions.

Definition 1 (convex functions). A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is **convex** on a set A if for all $x, y \in A$ and $\lambda \in [0, 1]$:

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$$

When this definition holds with the inequality being reversed, then f is said to be concave. From the definition, it is clear that a function f is convex if and only if $-f$ is concave.

(a) Why do we care whether a function is convex or not?

(b) Which of the following functions are convex? (Hint: draw a picture!)

(i) $|x|$

(ii) $\cos(x)$

(iii) $x^T x$

(c) Can a function be both convex and concave on the same set? If so, give an example. If not, describe why not.

4. Practical Methods for Checking Convexity

Using the definition to check whether a function is convex or not can be a tedious task in many situations. Some basic methods that can help us achieve the task in an efficient way are introduced below:

- for differentiable function, examine $f(y) \geq f(x) + \nabla f(x)^T(y - x)$ for any x, y in the domain of f .
- for twice differentiable functions, examine $\nabla^2 f(x) \succeq 0$ (i.e., the Hessian matrix is positive semidefinite).
- nonnegative weighted sum
- composition with affine function
- pointwise maximum and supremum

Note: there are even more such methods, which are covered in a convex optimization course or textbook.

(a) If f is differentiable, then f is convex if and only if $f(y) \geq f(x) + \nabla f(x)^T(y - x)$ for any x, y in the domain of f . A geometric interpretation of this characterization is that any tangent plane of a convex function f must lie entirely below f . One interesting application of this characterization is one of the most important inequalities in probability and statistics: the Jensen's inequality, which states that $\mathbb{E}f(X) \geq f(\mathbb{E}(X))$ when f is convex. Prove Jensen's inequality using the other inequality mentioned here.

(b) If f is twice differentiable with convex domain, then f is convex if and only if

$$\nabla^2 f(x) \succeq 0,$$

for any x in the domain of f . Use this method to show that the objective function in linear regression is convex.

(c) Let $\alpha \geq 0$ and $\beta \geq 0$, and if f and g are convex, then αf , $f + g$, $\alpha f + \beta g$ are all convex. One application: When a (possibly complicated) objective function can be expressed as a sum (e.g., the negative log-likelihood function), then showing the convexity of each individual term is typically easier.

(d) Suppose $f(\cdot)$ is convex, then $g(x) := f(Ax + b)$ is convex. Use this method to show that $\|Ax + b\|_1$ is convex (in x), where $\|z\|_1 = \sum_i |z_i|$.

(e) Suppose you know that f_1 and f_2 are convex functions on a set A . The function $g(x) := \max\{f_1(x), f_2(x)\}$ is also convex on A . In general, if $f(x, y)$ is convex in x for each y , then $g(x) := \sup_y f(x, y)$ is convex. Use this method to show that the largest eigenvalue of a matrix X , $\lambda_{\max}(X)$, is convex in X (Using the definition of convexity would make this question quite difficult).

(f) Does the same result hold for $h(x) := \min\{f_1(x), f_2(x)\}$? If so, give a proof. If not, provide convex functions f_1, f_2 such that h is not convex.

5. MAP as Regularization

Recall the regularization techniques that were presented in class this week and ponder their objectives:

(a) **Ridge-Regression:** $\hat{w}_{ridge} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2 + \lambda \|w\|_2^2$

(b) **LASSO:** $\hat{w}_{LASSO} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2 + \lambda \|w\|_1$

Reminder: don't ever regularize your bias term. This term doesn't add any complexity to the model (since it just shifts), so we'd like it to take on any value that best fits our training data.

The two types of regularization above can be derived from a statistical perspective in which we assume some prior belief about what the weights of our model should be and then observe data to further update the belief.

More specifically, let w denote our weights and Y, X our data (Y represents the labels and X the inputs). As before, $p(X, Y|w)$ represents the **likelihood function**. We specify our belief of what the weights should be through a **prior distribution** over $p(w)$. Using Bayes' Rule, we can write our updated belief of what the weights ought to be after observing the data as:

$$p(w|X, Y) = \frac{p(X, Y|w)p(w)}{p(X, Y)} = \frac{p(X, Y|w)p(w)}{\int_{w'} p(X, Y|w')p(w')dw'}$$

where we call $p(w|X, Y)$ the **posterior distribution** and $p(X, Y)$ the **evidence**.

What **Maximum A Posteriori Estimation (MAP)** does is compute the weights which maximize the posterior distribution, $p(w|X, Y)$. This type of estimation differs from MLE (which maximizes the likelihood function $p(X, Y|w)$) by taking into account our prior belief of what the weights are, namely $p(w)$. More specifically, the MAP estimate is:

$$\begin{aligned}\hat{w}_{MAP} &= \arg \max_w p(w|X, Y) \\ &= \arg \max_w \frac{p(X, Y|w)p(w)}{p(X, Y)} \\ &= \arg \max_w p(X, Y|w)p(w) \\ &= \arg \max_w \log p(X, Y|w) + \log p(w)\end{aligned}$$

where we dispose of the denominator because it doesn't depend on w . Contrast this with the MLE which is:

$$\hat{w}_{MLE} = \arg \max_w p(X, Y|w)$$

Let us now study how we can obtain the Ridge and LASSO regression objectives from this perspective:

(a) Suppose the elements of w are independently distributed according to a Laplacian distribution:

$$p(w_i) = \frac{\lambda}{4\sigma^2} \exp(-|w_i| \frac{\lambda}{2\sigma^2}).$$

Show that under this prior on w , MAP estimation of the linear measurement model recovers the LASSO objective.

(b) Derive an expression for the prior on w that corresponds to the ridge regression objective. What is the significance of this result?