

Definitions and notations

Zhehao Wang

November 1, 2018

1 Calculus and linear algebra

1.1 Limit

Let $f(x)$ be a function defined on an interval that contains $x = a$, except possibly at $x = a$, then we say that

$$\lim_{x \rightarrow a} f(x) = L$$

if for every $\epsilon > 0$ there is some number $\delta > 0$ such that

$$|f(x) - L| < \epsilon \text{ whenever } 0 < |x - a| < \delta$$

1.2 Gradient

Given $f(\vec{x})$ where $\vec{x} = (x_1, \dots, x_n)$ on \mathbb{R}^n

$$\nabla f(a_1, \dots, a_n) = \left(\frac{\partial f}{\partial x_1}(a_1, \dots, a_n), \dots, \frac{\partial f}{\partial x_n}(a_1, \dots, a_n) \right)$$

Intuition: gradient is a vector (the rate of change of your function, when you move in a certain direction), which in a two-dimensional space, tangents the curve at a given point.

1.3 Directional derivative

Homework notation

$$f'(x; u) = \lim_{h \rightarrow 0} \frac{f(x + hu) - f(x)}{h}$$

Wikipedias notation

$$\nabla_v f(\vec{x}) = \lim_{h \rightarrow 0} \frac{f(\vec{x} + h\vec{v}) - f(\vec{x})}{h}$$

Intuition: the rate-of-change of a function $f(\vec{x})$ on direction \vec{v} .

1.4 Vector norm

A norm is a function that assigns a strictly positive length or size to each vector in a vector space (except the zero vector which is assigned a length of 0).

Absolute value norm is a norm on the one-dimensional vector spaces formed by real or complex numbers.

$$\|x\| = |x|$$

Euclidean norm on a Euclidean space \mathbb{R}^n is such

$$\|\vec{x}\|_2 = \sqrt{x_1^2 + \dots + x_n^2}$$

Manhattan or taxicab norm

$$\|\vec{x}\|_1 = \sum_{i=1}^n |x_i|$$

p-norm

$$\|\vec{x}\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}}$$

Note that when $p = 1$, we get Manhattan norm, and when $p = 2$, we get Euclidean norm.

When $p = \infty$

$$\|\vec{x}\|_\infty = \max_i |x_i|$$

1.5 argmax

Points of the domain of some function at which the function values are maximized.

Given an arbitrary set X , a totally ordered set Y and a function $f : X \rightarrow Y$, the arg max over some subset S of X is defined by

$$\arg \max_{x \in S \subseteq X} f(x) = \{x \mid x \in S \wedge \forall y \in S : f(y) \leq f(x)\}$$

2 Probability

2.1 Random variable

A random variable $X : \Omega \rightarrow E$ is a measurable function from a set of possible outcomes Ω to a measurable space E . Often times $E = \mathbb{R}$

The probability that X takes on a value in a measurable set $S \subseteq E$ is written as

$$Pr(X \in S) = P(\omega \in \Omega \mid X(\omega) \in S)$$

Intuition: mapping outcomes of a random process to numbers, like this definition of X

$$X = \begin{cases} 0, & \text{if heads} \\ 1, & \text{if tails} \end{cases}$$

Instead of a traditional algebraic variable that can be solved for one value, a random variable can have different values (each with a probability) under different conditions.

2.2 Law of large numbers

X_1, X_2, \dots is an infinite sequence of independent and identically distributed (iid) random variables with expected value $E(X_1) = E(X_2) = \dots = \mu$, and

$$\bar{X}_n = \frac{1}{n}(X_1 + \dots + X_n) *$$

The weak law states that for any positive number ϵ

$$\lim_{n \rightarrow \infty} Pr(|\bar{X}_n - \mu| > \epsilon) = 0$$

The strong law states that

$$Pr(\lim_{n \rightarrow \infty} \bar{X}_n = \mu) = 1$$

Intuition: law of large numbers is a theorem that describes the average of the results obtained from a large number of trials should be close to the expected value, and will tend to become closer as more trials are performed.

3 Statistical learning

Let \mathbf{X} denote the input space, \mathbf{Y} denote the output space, and \mathbf{A} denote the action space.

A **decision function** / **prediction function** $f : \mathbf{X} \rightarrow \mathbf{A}$ maps an input to an action.

A **loss function** $l : \mathbf{A} \times \mathbf{Y} \rightarrow \mathbb{R}$ evaluates an action in the context of an output, and maps the error to a real. In traditional problems we can usually assume action has no effect on the output (like stock market prediction).

3.1 Risk, Bayesian function, empirical risk

Assume there is a data generating distribution $P_{\mathbf{X} \times \mathbf{Y}}$.

Risk of a decision function can then be defined as

$$R(f) = \mathbb{E}l(f(x), y)$$

Where an input/output pair (x, y) is generated independently and identically distributed from $P_{\mathbf{X} \times \mathbf{Y}}$.

*What does the sum of random variables mean? Seems that each X_i here means 'sampled' values following the distribution defined by the random variable, like, the outcome of an experiment

A **Bayesian decision function / target function** $f^* : \mathbf{X} \rightarrow \mathbf{A}$ is a function that achieves the minimal risk among all possible functions.

$$f^* = \arg \min_f R(f)$$

Risk cannot be calculated as we are not given $P_{\mathbf{X} \times \mathbf{Y}}$.
Let $\mathbf{D}_n = \{(x_1, y_1), \dots, (x_n, y_n)\}$ be drawn iid from $P_{\mathbf{X} \times \mathbf{Y}}$.
The **empirical risk** of $f : \mathbf{X} \rightarrow \mathbf{A}$ with respect to \mathbf{D}_n is

$$\hat{R}_n(f) = \frac{1}{n} \sum_{i=1}^n l(f(x_i), y_i)$$

By strong law of large numbers (2.2), when n is big enough, empirical risk can be used to approximate risk.

3.2 Empirical risk minimization, constrained empirical risk minimization

A function \hat{f} is an **empirical risk minimizer** if

$$\hat{f} = \arg \min_f \hat{R}_n(f)$$

The prediction function that produces the smallest empirical risk over set \mathbf{D}_n . This naturally leads to overfit.
So we minimize empirical risk, but only within a hypothesis space \mathbf{F} , being a set of prediction functions.

Constrained empirical risk minimizer can then be defined as

$$\hat{f}_n = \arg \min_{f \in \mathbf{F}} \frac{1}{n} \sum_{i=1}^n l(f(x_i), y_i)$$

If we plug Risk in instead, a **constrained risk minimizer** becomes

$$f_{\mathbf{F}}^* = \arg \min_{f \in \mathbf{F}} \mathbb{E}l(f(x), y)$$

3.3 Recap: linear regression with square loss