

Definitions and notations

Zhehao Wang

March 4, 2020

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 \mathbf{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

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

Intuition: the rate-of-change of a function $f(\vec{x})$ on direction \vec{v} . Remember that this is a scalar (since we are given the direction).

$$f'(x; u) = \nabla f(x)^T u$$

Meaning gradient on a certain direction vector u is $f(x)$'s directional derivative on u .

1.4 Taylor expansion

Every infinite differentiable function can be approximated by a series of polynomials.

Let $f(x)$ be a real function which is continuous on the closed interval $[a \dots b]$ and $n + 1$ times differentiable on the open interval $(a \dots b)$, let $\xi \in (a \dots b)$.

Then given any $x \in (a \dots b)$, there exists some $\eta \in \mathbb{R} : x \leq \eta \leq \xi$ or $\xi \leq \eta \leq x$ such that:

$$f(x) = \sum_{k=0}^n \frac{1}{k!} (x - \xi)^k f^{(k)}(\xi) + R_n$$

Where $R_n = \frac{1}{(n+1)!} (x - \xi)^{n+1} f^{(n+1)}(\eta)$ is the error term.

The expression

$$f(x) = \sum_{n=0}^{\infty} \frac{1}{n!} (x - \xi)^n f^{(n)}(\xi)$$

where n is taken to the limit is known as the Taylor series expansion of f about ξ .

1.5 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 \mathbf{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 $p = \infty$

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

1.6 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

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

1.7 Dot product, inner product

Maps two equal-size vectors to a real value.

$$\cdot : \mathbf{R}^d \times \mathbf{R}^d \rightarrow \mathbf{R}$$

Defined geometrically, given two vectors a and b ,

$$a \cdot b = \|a\| \|b\| \cos(\theta)$$

where θ is the angle between a and b .

Defined algebraically,

$$a \cdot b = a^T b$$

1.8 Inner product space / pre-Hilbert space, Hilbert space

An **inner product space** (over reals) is a vector space \mathbb{V} and an **inner product**, which is a mapping

$$\langle \cdot, \cdot \rangle : \mathbb{V} \times \mathbb{V} \rightarrow \mathbf{R}$$

This space has the properties that $\forall x, y, z \in \mathbb{V}$ and $a, b \in \mathbf{R}$,

- symmetry: $\langle x, y \rangle = \langle y, x \rangle$
- linearity: $\langle ax + by, z \rangle = a \langle x, z \rangle + b \langle y, z \rangle$
- positive-definiteness: $\langle x, x \rangle \geq 0$. And $\langle x, x \rangle = 0 \iff x = 0$

A **norm** on the inner product space can then be defined as

$$\|x\| = \sqrt{\langle x, x \rangle}$$

Parallelogram law states that a norm can be written in terms of an inner product on \mathbb{V} iff $\forall x, x' \in \mathbb{V}$,

$$2\|x\|^2 + 2\|x'\|^2 = \|x + x'\|^2 + \|x - x'\|^2$$

And if it can be, the inner product is given by the **polarization identity**

$$\langle x, x' \rangle = \frac{\|x\|^2 + \|x'\|^2 - \|x - x'\|^2}{2}$$

Pythagorean theorem states that two vectors are **orthogonal** if $\langle x, x' \rangle = 0$, denoted as $x \perp x'$. And x is orthogonal to a set S if for all $s \in S$, x is orthogonal to s .

If $x \perp x'$, then $\|x + x'\|^2 = \|x\|^2 + \|x'\|^2$

Roughly, **projection onto a plane** can then be defined as for $x \in \mathbb{V}$, let M be a subspace of inner product space \mathbb{V} , then m_0 is the projection of x onto M if $m_0 \in M$ and is the closest point to x in M . Meaning $\forall m \in M$, $\|x - m_0\| \leq \|x - m\|$.

A Hilbert space is a complete* inner product space. E.g. \mathbf{R}^d and the standard inner product. Any finite dimensional inner product space is a Hilbert space.

The **projection theorem** suggests, for a Hilbert space \mathbf{H} , M is a closed subspace of \mathbf{H} , for any $x \in \mathbf{H}$, there exists a unique $m_0 \in M$ for which

$$\|x - m_0\| \leq \|x - m\| \quad \forall m \in M$$

This m_0 is called the orthogonal projection of x onto M . Furthermore, $m_0 \in M$ is the projection of x onto M iff $x - m_0 \perp M$.

Projection reduces norm: let M be a closed subspace of \mathbf{H} . For any $x \in \mathbf{H}$, let m_0 be the projection of x onto M , then $\|m_0\| \leq \|x\|$ with equality only when $m_0 = x$.[†]

1.9 Cauchy-Schwarz inequality

Given two vectors u and v

$$|\langle u, v \rangle|^2 \leq |\langle u, u \rangle| \cdot |\langle v, v \rangle|$$

Where $|\langle u, u \rangle|$ is the inner product. The equality holds only when u and v are linearly independent (parallel).

1.10 Jacobian matrix

Jacobian matrix is the matrix of all first-order partial derivatives of a vector-valued function.

Suppose $\vec{f}: \mathbf{R}^n \rightarrow \mathbf{R}^m$, the Jacobian matrix \vec{J} of \vec{f} is defined as follows

$$\vec{J} = \begin{bmatrix} \frac{\partial \vec{f}}{\partial x_1} & \dots & \frac{\partial \vec{f}}{\partial x_n} \end{bmatrix}$$

or component-wise $\vec{J}_{ij} = \frac{\partial \vec{f}_i}{\partial x_j}$, meaning

*a space is complete if all Cauchy sequences in the space converge.

†think Pythagorean theorem on a \mathbf{R}^2 space

$$\vec{J} = \begin{bmatrix} \frac{\partial \vec{f}_1}{\partial x_1} & \cdots & \frac{\partial \vec{f}_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial \vec{f}_m}{\partial x_1} & \cdots & \frac{\partial \vec{f}_m}{\partial x_n} \end{bmatrix}$$

1.11 Order of a matrix

1.12 Inversibility of a matrix

1.13 Eigenvalue, eigenvector

An **eigenvalue** of a **linear transformation** * is a non-zero scalar that changes by only a scalar factor when that linear transformation is applied to it. Meaning the set of \vec{v} that satisfies

$$T(\vec{v}) = \lambda \vec{v}$$

Where T is the linear transformation, λ is a scalar and called **eigenvalue**. We can compute the eigenvalues λ of a square matrix A of size n by solving this **determinant**

$$|A - \lambda I| = 0$$

Where I is the **identity matrix (unit matrix)** of size n .

For each eigenvalue λ , we can then solve its (corresponding set of) eigenvectors \vec{v} using

$$(A - \lambda I) \cdot \vec{v} = 0$$

Intuition: applying a linear transformation on an eigenvector of that linear transformation yields a vector that is parallel to this eigenvector (multiplier: eigenvalue).

1.14 Positive-definite

A **symmetric** $n \times n$ real matrix is positive definite if the scalar $z^T M z$ is strictly positive for every non-zero column vector $z \in \mathbf{R}^n$.

The identity matrix I is positive definite.

For any real invertible matrix A , the product $A^T A$ is a positive definite matrix.

1.15 Convexity

A set C is **convex** if for any $x_1, x_2 \in C$ and any θ with $0 \leq \theta \leq 1$ we have

$$\theta x_1 + (1 - \theta) x_2 \in C$$

*Think square transformation matrix

Intuition: for all x_1, x_2 in C , all the points on the line segment connecting points x_1, x_2 are in C .

A **function** $f : \mathbf{R}^n \rightarrow \mathbf{R}$ is **convex** if $\text{dom } f$ is a convex set and if for all $x, y \in \text{dom } f$, and $0 \leq \theta \leq 1$, we have

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

Intuition: line segment connecting points x, y on the graph of f does not cross the graph of f .

Examples:

- $f(x) = ax + b$ is both convex and concave on \mathbf{R} for all $a, b \in \mathbf{R}$.
- $f(x) = |x|^p$ for $p \geq 1$ is convex on \mathbf{R} .
- $f(x) = e^{ax}$ for all a is convex on \mathbf{R} .
- Every norm on \mathbf{R}^n is convex.
- $f(x) = \max\{x\}$ is convex on \mathbf{R}^n .

A function f is **strictly convex** if the line segment connecting any two points on the graph of f lies strictly above the graph.

- When a function is convex, if there is a local minimum, then it is a global minimum.
- When a function is strictly convex, if there is a local minimum, then it is the unique global minimum.

For a multivariate twice differentiable function, $h : \Theta \subset \mathbb{R}^d \rightarrow \mathbb{R}$, $d \geq 2$, the gradient vector is

$$\nabla h(\theta) = \begin{pmatrix} \frac{\partial h}{\partial \theta_1}(\theta) \\ \vdots \\ \frac{\partial h}{\partial \theta_d}(\theta) \end{pmatrix}$$

A **Hessian matrix** of that function is

$$\nabla^2 h(\theta) = \begin{pmatrix} \frac{\partial^2 h}{\partial \theta_1 \partial \theta_1}(\theta) & \cdots & \frac{\partial^2 h}{\partial \theta_1 \partial \theta_d}(\theta) \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 h}{\partial \theta_d \partial \theta_1}(\theta) & \cdots & \frac{\partial^2 h}{\partial \theta_d \partial \theta_d}(\theta) \end{pmatrix}$$

h is concave $\implies x^T \nabla^2 h(\theta) x \leq 0 \ \forall x \in \mathbb{R}^d, \theta \in \Theta$.

h is strictly concave $\implies x^T \nabla^2 h(\theta) x < 0 \ \forall x \in \mathbb{R}^d, \theta \in \Theta$.

1.16 The general optimization problem

Standard form: minimize $f_0(x)$ subject to $f_i(x) \leq 0$, $i = 1, \dots, m$, $h_i(x) = 0$, $i = 1, \dots, m$.^{*} Where f_0 is the objective function and $x \in \mathbf{R}^n$ are the optimization variables.

We can replace $h(x) = 0$ with $h(x) \leq 0$ and $-h(x) \leq 0$, so we don't need the equality constraints.

The set of points satisfying the constraints is called the **feasible set**.

A point in the feasible set is called a **feasible point**.

If x is feasible and $f_i(x) = 0$, then we say inequality constraint $f_i(x) \leq 0$ is **active** at x .

The **optimal value** p^* of the problem is defined as

$$p^* = \inf \{f_0(x) | x \in \text{feasible set}\}$$

x^* is an **optimal point** (a solution to the problem) if x^* is feasible and $f(x^*) = p^*$.

1.17 Lagrangian duality

Given a general optimization problem:

$$\begin{aligned} &\text{minimize } f_0(x) \\ &\text{subject to } f_i(x) \leq 0, i = 1, \dots, m. \end{aligned}$$

The **Lagrangian** for it is defined as

$$L(x, \lambda) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x)$$

Where λ_i s are called **Lagrangian multipliers (dual variables)**.[†]

Supremum over Lagrangian gives back encoding of objective and constraints:

$$\begin{aligned} \sup_{\lambda \geq 0} L(x, \lambda) &= \sup_{\lambda \geq 0} (f_0(x) + \sum_{i=1}^m \lambda_i f_i(x)) \\ &= \begin{cases} f_0(x), & \text{when } f_i(x) \leq 0 \forall i \\ \infty, & \text{otherwise} \end{cases} \end{aligned}$$

And equivalent of the **primal form** of the optimization problem:

$$p^* = \inf_x \sup_{\lambda \geq 0} (f_0(x) + \sum_{i=1}^m \lambda_i f_i(x))$$

^{*} Assuming the intersections of domains of all f_i and h_i is not empty.

[†] Irrelevant with convexity.

We get the **Lagrangian dual problem** by swapping the inf and sup.

$$d^* = \sup_{\lambda \succeq 0} \inf_x (f_0(x) + \sum_{i=1}^m \lambda_i f_i(x))$$

Weak duality $p^* \geq d^*$ holds for any optimization problem. *

Duality gap is $p^* - d^*$, for convex problems, we often have **strong duality**: $p^* = d^*$.

The **Lagrangian dual function** is

$$g(\lambda) = \inf_x L(x, \lambda) = \inf_x (f_0(x) + \sum_{i=1}^m \lambda_i f_i(x))$$

The dual function is always concave.†

Lagrangian dual function gives a lower bound on optimal solution:

$$\begin{aligned} p^* &\geq d^* = \sup_{\lambda \succeq 0} g(\lambda) \\ p^* &\geq g(\lambda) \quad \forall \lambda \succeq 0 \end{aligned}$$

The **Lagrangian dual problem** is a search for best lower bound on p^* : maximize $g(\lambda)$ subject to $\lambda \succeq 0$.

λ is **dual feasible** if $\lambda \succeq 0$ and $g(\lambda) > -\infty$, and λ^* is **dual optimal** or **optimal Lagrange multipliers** if they are optimal for the Lagrange dual problem.

For a general optimization problem, if we have **strong duality**, we get a relationship called **complementary slackness** between

- the optimal Lagrange multiplier λ_i^* , and
- the i -th constraint at optimum: $f_i(x^*)$

$$\lambda_i^* f_i(x^*) = 0$$

Always have Lagrange multiplier being zero, or constraint is active at optimum, or both.

Lagrangian can be applied to illustrate the equivalence of Tikhonov and Ivanov forms of penalizing complexity measure.

1.18 Convex optimization standard form

Standard form of convex optimization problem: minimize $f_0(x)$ subject to $f_i(x) \leq 0, i = 1, \dots, m$. Where f_0, \dots, f_m are convex functions.

We usually have strong duality for convex optimization problems, but not always. The additional conditions needed are called **constraint qualifications**.

*Hint, for any general $f : W \times Z \rightarrow R$, this can be proved given $\inf_{w \in W} f(w, z_0) \leq f(w_0, z_0) \leq \sup_{z \in Z} f(w_0, z)$, add another sup to the leftmost term and inf to the rightmost term.

†Pointwise min of affine functions

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 $E = \mathbf{R}$

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

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

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 $\mathbb{E}[X_1] = \mathbb{E}[X_2] = \dots = \mu$ and $\mathbb{V}[X_1] = \dots = \sigma^2$, 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 weak law follows immediately from Chebychev's inequality.

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.

In other words, when n becomes large enough, average is a good i.e. consistent estimator of expectation.

*Or alternatively, $\lim_{n \rightarrow \infty} \frac{n(A)}{n} = P(A)$

2.3 Expectation, variance

Expectation

$$\begin{aligned}\mathbb{E}[X] &= \sum_{-\infty}^{+\infty} x P_X(x) \\ \mathbb{E}[\Phi(X)] &= \sum_{-\infty}^{+\infty} \Phi(x) P_X(x) \\ \mathbb{E}[\Phi(X, Y)] &= \sum_{-\infty}^{+\infty} \sum_{-\infty}^{+\infty} \Phi(x, y) P_{X,Y}(x, y)\end{aligned}$$

Where Φ is a function of random variables X and Y , $P_X(x)$ denotes $P(X = x)$, and $P_{X,Y}(x, y)$ is the joint distribution of X and Y ($P_{X,Y}(x, y) = P(X = x \text{ and } Y = y)$).

Chebyshev's inequality

We have $\mathbb{E}(X^2) = \sum_{-\infty}^{+\infty} x^2 P_X(x)$.

Given any number $\epsilon > 0$, construct X_1 from X s.t.

$$X_1 = \begin{cases} 0, & \text{if } |X| \leq \epsilon \\ \epsilon^2, & \text{if } |X| > \epsilon \end{cases}$$

Then obviously $X_1 \leq X^2$, and $\mathbb{E}[X_1] \leq \mathbb{E}[X^2]$, or equivalently

$$\epsilon^2 P(|X| > \epsilon) \leq \mathbb{E}[X^2]$$

Consequently

$$P(|X| > \epsilon) \leq \frac{1}{\epsilon^2} \mathbb{E}[X^2]$$

In particular if $\mathbb{E}[X^2] = 0$, then $P(|X| > \epsilon) = 0 \forall \epsilon > 0$ hence $X = 0$ with probability 1.

Variance

$$\mathbb{V}[X] = \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[x^2] - \mathbb{E}[x]^2$$

Derivation (recall that $\mathbb{E}[X]$ is a constant):

$$\begin{aligned}\mathbb{V}[X] &= \mathbb{E}[(X - \mathbb{E}[X])^2] \\ &= \mathbb{E}[X^2 - 2\mathbb{E}[X] \cdot X + \mathbb{E}[X]^2] \\ &= \mathbb{E}[X^2] - 2 \cdot \mathbb{E}[X] \cdot \mathbb{E}[X] + \mathbb{E}[X]^2 \\ &= \mathbb{E}[X^2] - \mathbb{E}[X]^2\end{aligned}$$

Some properties:

$$\mathbb{V}[X - a] = \mathbb{V}[X]$$

$$\mathbb{E}[cX] = c\mathbb{E}[X]$$

$$\mathbb{E}[X - a] = \mathbb{E}[X] - a$$

$$\mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y]$$

$$\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y]$$

if X and Y are independent.

$$\mathbb{V}[nX] = n^2\mathbb{V}[X]$$

$$\mathbb{V}[-X] = -\mathbb{V}[X]$$

$$\mathbb{V}[X + Y] = \mathbb{V}[X] + \mathbb{V}[Y]$$

If X and Y are independent.

It's then easy to show a "normalized" random variable $X' = \frac{X - \mathbb{E}[X]}{\sqrt{\mathbb{V}[X]}}$ satisfies $\mathbb{E}[X'] = 0$ and $\mathbb{V}[X'] = 1$.

Intuition of \mathbb{E} : say we are to use one value m to represent a random variable X , s.t. mean square error is minimized.

$$\begin{aligned} \text{MSE}(m) &= \mathbb{E}[(X - m)^2] \\ &= (\mathbb{E}[X - m])^2 + \mathbb{V}[X - m] \\ &= (\mathbb{E}[X] - m)^2 + \mathbb{V}[X] \end{aligned}$$

$$\frac{d\text{MSE}(m)}{dm} = 2m - 2\mathbb{E}[X]$$

Thus when $m = \mathbb{E}[X]$, $\frac{d\text{MSE}(m)}{dm} = 0$.

So the value μ we should use is the expectation.

The expectation can be estimated from the mean of samples (y_1, \dots, y_n) .

$$\hat{\mu} \equiv \frac{1}{n} \sum_{i=1}^n y_i$$

*

If samples are i.i.d, the **law of large number** suggests that

$$\hat{\mu} \rightarrow \mathbb{E}(Y) = \mu$$

Law of total expectation: for any random variables U and V ,

$$\mathbb{E}(U) = \mathbb{E}[\mathbb{E}[U|V]]$$

Now instead of predicting one value for a random variable, we are given input X and want to find a function $f : X \rightarrow Y$ s.t. the mean square error is minimal. This optimal function is the **regression function**, and

$$\begin{aligned} \text{MSE}(f) &= \mathbb{E}[(Y - f(X))^2] \\ &= \mathbb{E}[\mathbb{E}[(Y - f(X))^2|X]] \\ &= \mathbb{E}[\mathbb{V}[Y - f(X)|X] + (E[Y - f(X)|X])^2] \\ &= \mathbb{E}[\mathbb{V}[Y|X] + (E[Y - f(X)|X])^2] \end{aligned}$$

To minimize this, the first term in the expectation doesn't depend on our prediction, and the second term looks like previously when we are predicting the value of a random variable with one value, only with all expectations conditional on X .

So the optimal function $\mu(x) = \mathbb{E}[Y|X = x]$, and this is called **regression function**.

When Y depends on X (causal model), we can use a noise random variable ϵ (whose expectation is 0) and regression function $\mu(x)$ to estimate Y :

$$\mu(X) + \epsilon \rightarrow Y$$

When causal relationship $X \rightarrow Y$ cannot be assumed, the noise random variable can be dependent on X , thus the general form $Y|X = \mu(X) + \eta(X)$ (without loss of generality, $\eta(X)$'s expectation is 0).

2.4 Bias-variance trade-off

When X takes only a finite set of values and we have enough sample points, we can reliably approximate a regression function (law of large number).

But our sample points is almost always undersampled (e.g. when X is continuous), in which case we need interpolation, extrapolation and smoothing, and there are different methods to do these (different regression methods).

*Where the hat indicates empirical, as seen later on in empirical risk minimizer

Suppose the true regression function is $\mu(x)$ and we use $\hat{\mu}$ to make prediction.* The MSE of $\hat{\mu}$ (at x) can be written as such:

$$\begin{aligned}\text{MSE}(\hat{\mu}) &= \mathbb{E}[(Y - \hat{\mu}(x))^2] \\ &= \mathbb{E}[(Y - \mu(x) + \mu(x) - \hat{\mu}(x))^2] \\ &= \mathbb{E}[(Y - \mu(x))^2 + 2(Y - \mu(x))(\mu(x) - \hat{\mu}(x)) + (\mu(x) - \hat{\mu}(x))^2] \\ &= \mathbb{E}[\eta^2] + 2(\mu(x) - \hat{\mu}(x))\mathbb{E}[\eta] + \mathbb{E}[(\mu(x) - \hat{\mu}(x))^2] \\ &= \mathbb{V}[\eta] + (\mu(x) - \hat{\mu}(x))^2\end{aligned}$$

This is our first **bias-variance decomposition**, the second term is a bias by which our predictions are systematically off, and the first term is a variance that affect even the best prediction (which in general depends on x , let's denote it as σ_x^2).

In practice, $\hat{\mu}$ is not a single fixed function: it's something we estimate from sample data. If sample data are random, then exact regression function we get is random, too. Let's call this random function \hat{M}_n .[†]

Now if we are to analyze the prediction error of the *method*, averaging over all the possible training data sets

$$\begin{aligned}\text{MSE}(\hat{M}_n(x)) &= \mathbb{E}[(Y - \hat{M}_n(X))^2 | X = x] \\ &= \sigma_x^2 + (\mu(x) - \mathbb{E}[\hat{M}_n(x)])^2 + \mathbb{V}[\hat{M}_n(x)]\end{aligned}$$

This our second bias-variance decomposition. The first term is the same as before. The second term is in using \hat{M}_n to estimate μ , the **approximation error** / **approximation error**. The third term is the variance in our estimate of the regression function, even if we have an unbiased method $\mu(x) = \mathbb{E}[\hat{M}_n(x)]$, if there is a lot of variance in our estimates, we can expect to make large errors. **The catch** is, at least past a certain point, decreasing the approximation bias can only come through increasing the estimation variance. This is the **bias-variance trade-off**.

The trade off doesn't have to be one-for-one, sometimes we can lower the total error by introducing some bias as it gets rid of more variance than it adds approximation error.

In general, both approximation bias and estimation variance depend on n . A method is **consistent** when both of these go to 0 as $n \rightarrow \infty$. There can be multiple consistent methods for the same problem, and their biases and variances don't have to go to 0 at the same rate.

*As seen later, the true Bayesian form, and our prediction function which hopes to be as close as possible.

[†]The previous analysis really is $\text{MSE}(M_n(x) | M_n(x) = \hat{\mu})$

2.5 Ordinary least squares linear regression, linear smoothers

We choose to approximate $\mu(x)$ by $\alpha + \beta x$ and ask for the best a, b of the those constants.

For the sake of simplicity, we assume X is one dimensional and X and Y have mean 0.

$$\begin{aligned}\text{MSE}(\alpha, \beta) &= \mathbb{E}[(Y - \alpha - \beta X)^2] \\ &= \mathbb{E}[\mathbb{V}[Y|X]] + \mathbb{E}[(\mathbb{E}[Y - \alpha - \beta X|X])^2]\end{aligned}$$

The first term doesn't depend on α or β , so we can drop it. Taking the derivative of the second term, we have

$$\frac{\partial \text{MSE}}{\partial \alpha} = \mathbb{E}[-2 \cdot (Y - \alpha - \beta X)]$$

and $a = \mathbb{E}[Y] - b\mathbb{E}[X] = 0$. (we assumed X and Y both have mean 0)

$$\frac{\partial \text{MSE}}{\partial \beta} = \mathbb{E}[-2X \cdot (Y - \alpha - \beta X)]$$

and $b = \frac{\mathbb{E}[XY]}{\mathbb{E}[X^2]} = \frac{\text{Cov}[X,Y]}{\mathbb{V}[X]} = \frac{\sum_i y_i x_i}{\sum_i x_i^2}$. (we assumed X and Y both have mean 0)

And we are now in a position to see how the least square linear regression model is really a smoothing of the data:

$$\hat{\mu}(x) = \hat{b}x = \sum_i y_i \frac{x_i}{\sum_j x_j^2} x = \sum_i y_i \frac{x_i}{n\hat{\sigma}_X^2} x$$

Where $\hat{\sigma}_X^2$ is the sample variance of X .

The **intuition** is that our prediction is a weighted average of observed values y_i of the dependent variable, where the weights are proportional to how far x_i is from the center (relative to the variance), and proportional to the magnitude of x .

The linear regression line is a special case of **linear smoothers**, which are estimates of the regression function with the following form:

$$\hat{\mu}(x) = \sum_i y_i \hat{w}(x_i, x)$$

They are called linear as the predictions are linear in the responses y_i , as functions of x they are generally nonlinear. In linear regression's case, as shown earlier,

$$\hat{w}(x_i, x) = \frac{x_i}{n\hat{\sigma}_X^2} x$$

This ignores the distance between x_i and x .

2.6 k-Nearest-Neighbor Regression

Consider **nearest-neighbor regression**, which is very sensitive to the distance between x_i and x :

$$\hat{w}(x_i, x) = \begin{cases} 1 & x_i \text{ nearest neighbor of } x \\ 0 & \text{otherwise} \end{cases}$$

This will include the noise into its prediction. We might instead do k -nearest neighbor regression (as noise tend to cancel each other out), where as we increase k we get smoother functions, until $k = n$ where we get a constant.

$$\hat{w}(x_i, x) = \begin{cases} \frac{1}{k} & x_i \text{ one of the } k \text{ nearest neighbor of } x \\ 0 & \text{otherwise} \end{cases}$$

Because k -nearest-neighbors averages over only a fixed number of neighbors, each of which is a noisy sample, it always has some noise in its prediction, and is generally not consistent.

2.7 Kernel smoothers

With kNN regression each testing point is predicted using information from only a few data points, say we want to use all the training data like ordinary linear regression, but in a location-sensitive way.

Kernel smoothing does this. We pick a kernel function $K(x, x)$ which satisfies

- $K(x_i, x) \geq 0$
- $K(x_i, x)$ depends only on the distance $x_i - x$, not the individual arguments. (Hence we can write K as a one-argument function, $K(x_i - x)$)
- $\int x K(0, x) dx = 0$ and
- $0 < \int x^2 K(0, x) dx < \infty$

These conditions together imply $|x_i - x| \rightarrow \infty, K(x_i, x) \rightarrow 0$.*

The Nadaraya-Watson estimate of the regression function is

$$\hat{\mu}(x) = \sum_i y_i \frac{K(x_i, x)}{\sum_j K(x_j, x)}$$

$K(x_i, x)$ is large if x_i is close to x , so this places a lot of weight on training data points close to the point where we are trying to predict.

*Example of such functions include the density of the uniform distribution $(-\frac{h}{2}, \frac{h}{2})$, and the density of the standard Gaussian $\mathbf{N}(0, \sqrt{h})$ distribution. h can be any positive number and is called **bandwidth**, which controls the degree of smoothing, $h \rightarrow \infty$ we revert to taking the global mean, and $h \rightarrow 0$ results in spikier functions.

2.8 Bernoulli trials, Binomial and Poisson distributions, Normal distribution

A **Bernoulli trial** means identical independent experiments in each of which an event A may occur with probability $P(A)$. In the case of n consecutive Bernoulli trials, each elementary event ω can be described by a sequence of n 0s and 1s. Because of the independence of the trials, the probability of such an event ω with k successes and $n - k$ failures is $p^k(1 - p)^{n-k}$ where p is the probability of event A in one trial.

Let random variable ξ denote the number of successes in n Bernoulli trials, $\xi(\omega) = k$ if k successes occur in the event ω .

$$P(\xi = k) = \binom{n}{k} p^k (1 - p)^{n-k} = \frac{n!}{k!(n-k)!} p^k (1 - p)^{n-k}$$

This is known as the **binomial distribution**, and specified by two parameters n and p .

The random variable ξ can be seen as the sum of n independent random variables $\xi = \xi_1 + \dots + \xi_n$, where each ξ_i has $P(\xi_i = 1) = p$ and $P(\xi_i = 0) = 1 - p$. It's easy to show that $\mathbb{E}[\xi_i] = p$ and $\mathbb{V}[\xi_i] = p(1 - p)$, and $\mathbb{E}[\xi] = np$, $\mathbb{V}[\xi] = np(1 - p)$. Approximation when n is large and p is small.

$$p_\xi(k) \approx \frac{(np)^k}{k!} e^{-np}$$

A random variable ξ taking integral values $0, 1, 2, \dots$ is said to have **Poisson distribution** if

$$p_\xi(k) \approx \frac{a^k}{k!} e^{-a}$$

The distribution is specified by a single positive parameter $a = \mathbb{E}[\xi]$. Poisson distribution is the same as the approximation of binomial distribution when the number of trials is large and probability of success is small.

A **normal distribution** probability density follows the form

$$p(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

And distribution function follows the form

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{x^2}{2}} dx$$

This has expectation 0 and variance 1.

More generally, a normal random variable is a random variable with probability density

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-a)^2}{2\sigma^2}}$$

It has expectation a and variance σ^2 .

2.9 Uncorrelated, independent, Eve's law, Central Limit Theorem

Random variables X and Y are **independent**:

$$F(X, Y) = F(X)F(Y)$$

Where $F(X)$, $F(Y)$ are cumulative distribution functions of X , Y ; and $F(X, Y)$ is the joint cumulative distribution function,

$$F(X, Y) = P(X \leq x, Y \leq y)$$

Random variables X and Y are **uncorrelated**:

$$Cov(X, Y) = 0$$

Where the covariance $Cov(X, Y) = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$

Eve's law

$$\mathbb{V}[Y] = \mathbb{E}[\mathbb{V}[Y|X]] + \mathbb{V}[\mathbb{E}[Y|X]]$$

Correspondingly, some call the law of total expectation **Adam's law**.

2.10 Generating function

Let ξ be a discrete random variable taking values $0, 1, 2, \dots$ with probabilities $P_\xi(k) = Pr(\xi = k)$, $k = 0, 1, 2, \dots$, then the function

$$F_\xi(z) = \sum_{k=0}^{\infty} P_\xi(k) z^k, \quad |z| \leq 1$$

is called the **generating function** of the random variable ξ .

The probability distribution of the random variable ξ is uniquely determined by its generating function $F_\xi(z)$:

$$P_\xi(k) = \frac{1}{k!} F_\xi^{(k)}(0), \quad k = 0, 1, 2, \dots$$

Where $F_\xi^{(k)}(0)$ is the k -th derivative of $F_\xi(z)$.

Since $\mathbb{E}[\Phi(x)] = \sum_{k=-\infty}^{+\infty} \Phi(k)P(k)$, $F_\xi(z)$ at a fixed z is the expectation of the random variable $\Phi(\xi) = z^\xi$:

$$F_\xi(z) = \mathbb{E}[z^\xi], \quad |z| \leq 1$$

The sequence of probability distribution $P_n(k)$, $n = 1, 2, \dots$ with generating functions $F_n(z)$, $n = 1, 2, \dots$ converges weakly to the limiting distribution $P(k)$ iff $\lim_{n \rightarrow \infty} F_n(z) = F(z)$ where $F(z)$ is the generating function of $P(k)$.

Given a real random variable ξ , **characteristic function** of ξ is the function

$$f_\xi(t) = \mathbb{E}[e^{i\xi t}], \quad -\infty < t < \infty$$

For a discrete random variable the characteristic function $f_\xi(z)$ on the boundary of the unit circle $|z| = 1$, i.e.

$$f_\xi(t) = F_\xi(e^{it}) = \sum_{k=0}^{\infty} P_\xi(k) e^{ikt}$$

This represents $f_\xi(t)$ as Fourier series with probability $P(\xi = k)$ as its coefficients.

For a continuous random variable, the characteristic function is the **Fourier transform** of the density $p_\xi(x)$:

$$f_\xi(t) = \int_{-\infty}^{+\infty} P_\xi(x) e^{ixt} dx$$

In both cases, $P_\xi(k)$ is uniquely determined by the characteristic function.

2.11 Central limit theorem

Given a sequence of i.i.d r.v. $X_i, i = 1, 2, \dots, n$ with mean $\mu = \mathbb{E}[X]$ and variance $\sigma^2 = \mathbb{V}[X]$, the weak and strong laws of large number says

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow[n \rightarrow \infty]{\mathbb{P}, a.s.} \mu$$

And central limit theorem says

$$\sqrt{n} \frac{\bar{X}_n - \mu}{\sigma} \xrightarrow[n \rightarrow \infty]{(d)} N(0, 1)$$

Intuition: LLN states that average is a consistent estimator of expectation when n is large enough. CLT says how good it is, and that the distribution of the sum of a large number of iid random variables is approximately normal.

Suppose the sequence $\xi_k, k = 1, 2, \dots$ with finite means $a_k = \mathbb{E}[\xi_k]$ and variance $\sigma_k^2 = \mathbb{V}[\xi_k]$ satisfies **Lyapunov condition**

$$\lim_{n \rightarrow \infty} \frac{1}{B_n^2} \sum_{k=1}^n \mathbb{E}[(\xi_k - a_k)^2] = 0$$

where $B_n^2 = \mathbb{V}[S_n] = \sum_{k=1}^n \sigma_k^2$

Then the sequence of random variables satisfy the central limit theorem.

2.12 Hoeffding's inequality

Let n be a positive integer and X, X_1, \dots, X_n be i.i.d r.v. s.t. $X \in [a, b]$. Let $\mu = E[X]$. Then for all $\epsilon > 0$,

$$P(|\bar{X}_n - \mu| \geq \epsilon) \leq 2e^{-\frac{2n\epsilon^2}{(b-a)^2}}$$

Intuition: average of bounded i.i.d r.v.'s probability distribution function is bounded by a Gaussian distribution (thinner tails) (even without n approaching infinity).

2.13 Different types of convergence

Let $(T_n)_{n \geq 1}$ be a sequence of r.v. and T a r.v. (T can be deterministic)

- Almost surely (a.s.) convergence: $T_n \xrightarrow[n \rightarrow \infty]{a.s.} T$ iff $\mathbf{P}[\{\omega : T_n(\omega) \xrightarrow[n \rightarrow \infty]{} T(\omega)\}] = 1$
- Convergence in probability: $T_n \xrightarrow[n \rightarrow \infty]{\mathbf{P}} T$ iff $\mathbf{P}[|T_n - T| \geq \epsilon] \xrightarrow[n \rightarrow \infty]{} 0, \forall \epsilon > 0$
- Convergence in L^p : $T_n \xrightarrow[n \rightarrow \infty]{(L^p)} T$ iff $\mathbf{P}[|T_n - T|^p] \xrightarrow[n \rightarrow \infty]{} 0$
- Convergence in distribution L^p : $T_n \xrightarrow[n \rightarrow \infty]{(d)} T$ iff $\mathbf{P}[T_n \leq x] \xrightarrow[n \rightarrow \infty]{} \mathbf{P}[T \leq x]$

Intuition: the first is very strong. It says if you measure one under whatever observation (ω), it's almost the same as if you measure the other. One implies two. For three if you converge in L^p you also converge in L^q for all $q < p$. Two implies four: if you converge in probability, you converge in distribution.

Continuous mapping theorem: if f is a continuous function, then

$$T_n \xrightarrow[n \rightarrow \infty]{a.s./\mathbf{P}/(d)} T \implies f(T_n) \xrightarrow[n \rightarrow \infty]{a.s./\mathbf{P}/(d)} f(T)$$

Intuition: apply a continuous function on the sequence of r.v. and convergence still holds.

One can add, multiply, divide, limit, etc, when converge almost surely and in probability. If $U_n \xrightarrow[n \rightarrow \infty]{a.s./\mathbf{P}} U$ and $V_n \xrightarrow[n \rightarrow \infty]{a.s./\mathbf{P}} V$, then

$$V_n + U_n \xrightarrow[n \rightarrow \infty]{a.s./\mathbf{P}} V + U$$

$$V_n U_n \xrightarrow[n \rightarrow \infty]{a.s./\mathbf{P}} V U$$

Same goes for division if $V \neq 0$. Strong LLN suggests almost surely convergence, and weak LLN suggests convergence in probability.

Slutsky's theorem suggests under certain scenario, you can combine two sequences when one converges in probability and the other converges in distribution.

By Slutsky, we have the **Delta method**, which allows us to apply functions on sequence of r.v. in CLT. I.e.

Let $(Z_n)_{n \geq 1}$ be a sequence of r.v. that asymptotically normal around θ , i.e. satisfies

$$\sqrt{n}(Z_n - \theta) \xrightarrow[n \rightarrow \infty]{a.s./\mathbf{P}} \mathbf{N}(0, \sigma^2)$$

Let $g : \mathbb{R} \rightarrow \mathbb{R}$ be continuously differentiable at the point θ , then $g(Z_n)$ is also asymptotically normal, more precisely

$$\sqrt{n}(g(Z_n) - g(\theta)) \xrightarrow[n \rightarrow \infty]{a.s./\mathbf{P}} \mathbf{N}(0, g'(\theta)^2 \sigma^2)$$

2.14 Degree of freedom

The number of values in the final calculation of a statistic that are free to vary without violating constraints.

In general, the **Degrees of Freedom** of an estimate of a parameter are equal to the number of independent score that go into the estimate minus the number of parameters used as intermediate steps in the estimation of the parameter itself. Hence most of the time sample variance has $N - 1$ degrees of freedom, N random points minus the 1 intermediate step, the sample mean.

2.15 Statistical experiment

The ultimate goal of statistics is to say what distribution your data comes from. The purpose of modeling is to restrict the space of possible distributions to a subspace that is possible to estimate.

In that sense, all models are wrong, some models are useful.

A **statistical experiment** is the pair $(E, \{\mathbf{P}_\theta\}_{\theta \in \Theta})$, where E is a set, the **sample space** (the possible values my observations can take), and each of P_θ is a probability distribution (e.g. $Poiss(\theta)$, $Bernoulli(\theta)$, etc. θ is a parameter set, it can be a pair, a couple of numbers, etc. The point is, once θ is known, you know the distribution perfectly.)

Well specified means for an observation X , $\exists \theta \in \Theta$ where $X \sim \mathbf{P}(\theta)$. This is a strong assumption. (X is just any one among the i.i.d $X_1, \dots X_n$).

θ (θ^*) is the true parameter we want to estimate. θ may have physical links to the phenomenon, like in linear regression, it suggests how sensitive it is to a dimension.

A **parametric model / statistics** is where θ has finite number of values θ_d , and all $\theta_d \in \mathbb{R}$.

Non-parametric estimation can be, e.g. when you have infinite numbers in θ to

estimate, or when θ is a function.

Example models:

- A Bernoulli trial is a pair $(\{0, 1\}, \{Ber(p)\}_{p \in (0,1)})$.
- An exponential trial is a pair $((0, +\infty], \{Exp(\lambda)\}_{\lambda > 0})$.
- A Poisson trial is a pair $(\mathbb{N}, \{Poiss(\lambda)\}_{\lambda > 0})$.
- A normal distribution is a pair $(\mathbb{R}, \{\mathbf{N}(\mu, \sigma^2)\}_{\mu \in \mathbb{R}, \sigma^2 > 0})$
- A Gaussian with known distribution (imagine machine parts with manufacturer labeled standard deviation), $(\mathbb{R}, \mathbf{N}(\mu, \sigma^2)_{\mu \in \mathbb{R}})$, compared with previous, the parameter space only need to specify μ , as σ^2 is known
- Censored data (imagine asking people's salary and getting answers like 5 figures). Given an indicator function I ($X = I(Y > 5)$, meaning $X = 1$ if $Y > 5$ else $X = 0$), and $Y \sim Exp(\lambda)$, then we have $(\{0, 1\}, \{Ber(e^{-5\lambda})\}_{\lambda > 0})$, the parameter space can also be expressed in p , the λ comes from integral of exponential on $(5, \lambda)$.
- A uniform distribution $X \sim \mathbf{U}([0, \theta])$, $\theta > 0$ is a pair $([0, +\infty), \{Unf(\theta)\}_{\theta > 0})$. Note that the sample space is not $[0, \theta]$ as it should not depend on θ .
- Or an arbitrary pmf (probability mass function) for a sample space of $1 \dots 100$.

The parameter θ is called **identified / identifiable** (in this model) iff the map $\theta \in \Theta \rightarrow \mathbf{P}_\theta$ is **injective**. I.e.

$$\theta \neq \theta' \implies \mathbf{P}_\theta \neq \mathbf{P}_{\theta'}$$

Intuition: it could be the case that two different θ s can give me the exact same probabilities, in which case θ is not identifiable (even with infinite amount of data). In the above examples the parameters are identifiable.

($\mathbf{N}(\mu, \sigma^2)$ minus μ then divide by σ gives you standard Gaussian.)

Given a statistical model, we want to estimate θ .

An **estimator** is a measurable function of data $\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$ a.k.a statistic. (e.g. average, max, etc. Measurable means that given data, you can compute it, it does not depend on θ . Things that are implicitly defined, such as sup or inf, are not measurable.)

An estimator is **weakly consistent** if it converges in probability to θ , or strongly consistent if it converges almost surely. (When you get more data, the estimator converges to the true parameters.)

The **bias** of an estimator is how far the $\hat{\theta}$ it gets is from θ : $\mathbb{E}[\hat{\theta}] - \theta$. If the above is zero, estimator is unbiased. (If I were to repeat the experiment over and over again, on average we get the true parameter. But often times you don't get to repeat the experiment.)

However being unbiased is not the only thing we are after: consider a Bernoulli distribution, for n i.i.d data points, by linearity $\mathbb{E}[\bar{X}_n] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_i] = \mathbb{E}[X_1] = p$. This is saying the bias of n data points is as good as 1 data point.

Hence we also have **(quadratic) risk** of $\hat{\theta}$, $\mathbb{E}[|\hat{\theta} - \theta|^2]$. If this goes to 0 as $n \rightarrow \infty$, then this is saying $\hat{\theta} \xrightarrow[n \rightarrow \infty]{L^2} \theta$, and this implies convergence in probability,

hence $\hat{\theta}$ is weakly consistent. So going back to the Bernoulli example, $\mathbb{E}[\bar{X}_n]$ is consistent as $n \rightarrow \infty$, but $\mathbb{E}[X_1]$ is not.

This risk is the sum of bias and variance.

$$\begin{aligned} \mathbb{E}[|\hat{\theta} - \theta|^2] &= \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}] + \mathbb{E}[\hat{\theta}] - \theta)^2] \\ &= \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2] + (\mathbb{E}[\hat{\theta}] - \theta)^2 \\ &= \mathbf{V}(\hat{\theta}) + \text{bias}^2 \end{aligned}$$

E.g. given $X_1 \dots X_n \sim \text{Ber}(\theta)$,

- the estimator \bar{X}_n has 0 bias and $\frac{1}{n}\theta(1-\theta)$ variance. Its risk is the sum of the two.
- the estimator 0.5 has bias $(0.5 - \theta)$ and variance 0.
- the estimator X_1 has bias 0 and variance $\theta(1-\theta)$. In this sense the first estimator has smaller risk.

Bias variance tradeoff happens all the time especially in high dimensional and non-parametric estimations.

2.16 Confidence interval

Let $(E, \{\mathbb{P}_\theta\}_{\theta \in \Theta})$ be a statistical model based on observations X_1, \dots, X_n and assume $\Theta \subseteq \mathbb{R}$.

Let $\alpha \in (0, 1)$, **confidence interval** of level $1 - \alpha$ for θ means any random (i.e. depending on $X_1 \dots X_n$) interval I whose boundaries do not depend on θ and such that

$$\mathbb{P}_\theta[I \ni \theta] \geq 1 - \alpha, \quad \forall \theta \in \Theta$$

Confidence interval of asymptotic level $(1 - \alpha)$ for θ means any random (i.e. depending on $X_1 \dots X_n$) interval I whose boundaries do not depend on θ and such that

$$\lim_{n \rightarrow \infty} \mathbb{P}_\theta[I \ni \theta] \geq 1 - \alpha, \quad \forall \theta \in \Theta$$

(when you use CLT.)

Intuition: you are given samples and want to understand something about the underlying distribution, in particular, you have sample mean and want to understand something about population mean. Given a level of confidence $x\%$, you want to find the interval which contains the population mean with $x\%$.

When dealing with normal distribution, confidence interval depends on whether population standard deviation is known or estimated with sample standard deviation. If not, then student's T distribution is used as the critical value.

2.17 Maximum likelihood estimation

Let $(E, \{\mathbb{P}_\theta\}_{\theta \in \Theta})$ be associated with i.i.d r.v. X_1, \dots, X_n and assume $\exists \theta^* \in \Theta$ such that $X_i \sim \mathbb{P}_{\theta^*}$: θ^* is the true parameter.

Statistician's goal is, given X_1, \dots, X_n , find an estimator $\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$ such that $\mathbb{P}_{\hat{\theta}}$ is close to \mathbb{P}_{θ^*} . This means $|\mathbb{P}_{\hat{\theta}}(A) - \mathbb{P}_{\theta^*}(A)|$ is small for all $A \subset E$.

The total variation distance is then defined as $TV(\mathbb{P}_\theta, \mathbb{P}_{\theta'}) = \max |\mathbb{P}_\theta(A) - \mathbb{P}_{\theta'}(A)| \quad \forall A \subset E$

This can then be translated to sum on the two pmf for discrete r.v., or integral on two density functions for continuous r.v..

Given the PMF of X : $\mathbb{P}_\theta(X = x) = p_\theta(x)$ for all $x \in E$,

$$p_\theta(x) \geq 0, \quad \sum_{x \in E} p_\theta(x) = 1$$

The total variation distance between \mathbb{P}_θ and $\mathbb{P}_{\theta'}$ is a function of PMF's p_θ and $p_{\theta'}$:

$$TV(\mathbb{P}_\theta, \mathbb{P}_{\theta'}) = \frac{1}{2} \sum_{x \in E} (p_\theta(x) - p_{\theta'}(x))$$

And in the continuous case f_θ being the density function,

$$TV(\mathbb{P}_\theta, \mathbb{P}_{\theta'}) = \frac{1}{2} \int_E |f_\theta(x) - f_{\theta'}(x)| dx$$

Properties of total variation distance include: symmetric, non-negative, definite ($TV(\mathbb{P}_\theta, \mathbb{P}_{\theta'}) = 0 \implies \mathbb{P}_\theta = \mathbb{P}_{\theta'}$), and triangle inequality $TV(\mathbb{P}_\theta, \mathbb{P}_{\theta'}) \leq TV(\mathbb{P}_\theta, \mathbb{P}_{\theta''}) + TV(\mathbb{P}_{\theta'}, \mathbb{P}_{\theta''})$

These implies total variation is a **distance**. We then want to minimize the total variation distance of θ and θ^* , but we don't know θ^* . We build an estimator $\hat{TV}(\mathbb{P}_\theta, \mathbb{P}_{\theta^*})$ for all $\theta \in \Theta$. Then find $\hat{\theta}$ that minimizes function $\theta \rightarrow \hat{TV}(\mathbb{P}_\theta, \mathbb{P}_{\theta^*})$.

The problem then becomes building the above estimator.

Kullback-Leibler (KL) divergence between two distributions is defined by the following. Discrete:

$$KL(\mathbb{P}_\theta, \mathbb{P}_{\theta'}) = \sum_{x \in E} p_\theta(x) \log\left(\frac{p_\theta(x)}{p_{\theta'}(x)}\right)$$

Or continuous:

$$KL(\mathbb{P}_\theta, \mathbb{P}_{\theta'}) = \int_E f_\theta(x) \log\left(\frac{f_\theta(x)}{f_{\theta'}(x)}\right) dx$$

KL makes the distance (e.g. between two Gaussians) easy to compute. KL is not symmetric, and it is not a distance, but definite property still holds. KL looks just like the density of the function log division. And,

$$KL(\mathbb{P}_\theta, \mathbb{P}_{\theta'}) = \mathbb{E}_{\theta*}[\log(\frac{p_{\theta*}(X)}{p_\theta(X)})] = \mathbb{E}_{\theta*}[\log(p_{\theta*}(X))] - \mathbb{E}_{\theta*}[\log(p_\theta(X))]$$

The first term is constant (to θ change), and the second term can be estimated with LLN, and we get

$$\hat{KL}(\mathbb{P}_{\theta*}, \mathbb{P}_\theta) = \text{constant} - \frac{1}{n} \sum_{i=1}^n \log(p_\theta(X_i))$$

We still don't know the constant / first term yet, but the minimizer θ of $\hat{KL}(\mathbb{P}_{\theta*}, \mathbb{P}_\theta)$ is the same no matter the constant.

$$\arg \min_{\theta \in \Theta} \hat{KL}(\mathbb{P}_{\theta*}, \mathbb{P}_\theta) \iff \arg \max_{\theta \in \Theta} \prod_{i=1}^n p_\theta(X_i)$$

This is the **maximum likelihood principle**.

Intuition: the multiplicative form we arrived at is the joint density of $X_1 \dots X_n$. Imagine the model is a Gaussian centered at θ^* , and we have only one observation at θ' , then the maximum likelihood estimator is going to estimate a Gaussian centered around θ' . Then if we have n points, we are just multiplying the values mapped onto the vertical line of different θ s, and finding the θ such that the product is maximized.

In the real world maximum likelihood estimators often don't have a closed form, hence the many optimizations in machine learning.

The **likelihood of a model** $(E, (\mathbb{P}_\theta)_{\theta \in \Theta})$ associated with a sample of i.i.d r.v X_1, \dots, X_n can then be defined as

$$\begin{aligned} L_n : E^n \times \Theta &\rightarrow \mathbb{R} \\ (x_1, \dots, x_n, \theta) &\rightarrow \mathbb{P}_\theta[X_1 = x_1, \dots, X_n = x_n] \end{aligned}$$

If E is discrete.

$$\begin{aligned} L_n : E^n \times \Theta &\rightarrow \mathbb{R} \\ (x_1, \dots, x_n, \theta) &\rightarrow \prod_{i=1}^n f_\theta(x_i) \end{aligned}$$

where f is the density. If E is continuous. This can be derived from KL divergence.

Example: likelihood of a Bernoulli model (discrete), Gaussian (continuous)

In stats, you first design your estimator as a function of random variables, then you get data and plug them in. (Keep the random variables for as long as you can.)

The **maximum likelihood estimator** of θ is defined as

$$\hat{\theta}_n^{MLE} = \arg \max_{\theta \in \Theta} L(X_1, \dots, X_n, \theta)$$

In practice, we use the log likelihood estimator as argmax is the same

$$\hat{\theta}_n^{MLE} = \arg \max_{\theta \in \Theta} \log L(X_1, \dots, X_n, \theta)$$

Maximum likelihood estimation is what's driving our intuition, e.g. when we use the average to estimate the p of a Bernoulli trial.

MLE minimizes KL divergence.

The **Fisher information** of a statistical model $l(\theta) = \log L_1(X, \theta)$ is defined as

$$\begin{aligned} I(\theta) &= \mathbb{E}[\nabla l(\theta) \nabla l(\theta)^T] - \mathbb{E}[\nabla l(\theta)] \mathbb{E}[\nabla l(\theta)]^T \\ &= -\mathbb{E}[\nabla^2 l(\theta)] \end{aligned}$$

Assuming l is twice differentiable.

If $\theta \in \mathbb{R}$, we get:

$$I(\theta) = \text{var}[l'(\theta)] = -\mathbb{E}[l''(\theta)] = \int \frac{(\frac{\partial f_\theta(x)}{\partial \theta})^2}{f_\theta(x)} dx$$

Intuition what Fisher information says is the more curvy the distribution is, the flatter it is when your estimator is around maximum, and the less sensitive we are to fluctuations (and the larger Fisher information, i.e. how much information of θ your model has).

Let $\theta^* \in \Theta$, assume the model is identified, $\forall \theta \in \Theta$, the support of \mathbb{P}_θ does not depend on θ , θ^* is not on the boundary of Θ , $I(\theta)$ is invertible in a neighborhood of θ^* , then $\hat{\theta}_n^{MLE}$ satisfies

$$\hat{\theta}_n^{MLE} \xrightarrow[n \rightarrow \infty]{\mathbb{P}} \theta^* \text{ w.r.t } \mathbb{P}_{\theta^*}$$

MLE converges in probability to θ^* . Bias goes to 0.

$$\sqrt{n}(\hat{\theta}_n^{MLE} - \theta^*) \xrightarrow[n \rightarrow \infty]{(d)} N(0, I(\theta^*)^{-1}) \text{ w.r.t } \mathbb{P}_{\theta^*}$$

MLE minus θ^* converges in distribution to normal distribution. θ^* may not even have a closed form, and we still have the same properties of an average. CLT-like. Variance is measured by Fisher information.

2.18 Method of moments

2.19 Markov chains

Consider a physical system with the following properties

- Can occupy any of a finite or countably infinite number of states $\epsilon_1, \epsilon_2, \dots$
- Starting from some initial state at time $t = 0$, change its state randomly at the times $t = 1, 2, \dots$. Let $\xi(t)$ denote the state of the system at time t , the evolution of the system can be described in steps $\xi(0) \rightarrow \xi(1) \dots$
- At $t = 0$, the system occupies the state ϵ_i with initial probability $p_i^0 = P(\xi(0) = \epsilon_i), i = 1, 2, \dots$
- Suppose at $t = n$ the system is at state ϵ_i , then the probability that the system goes into state ϵ_j at the next step is given by transition probability $p_{ij} = P(\xi(n+1) = \epsilon_j | \xi(n) = \epsilon_i)$, regardless of its behavior before time n .

A random process described by this model is a **Markov chain**.

p_{ij} can define a transition matrix P . Given $P(0) = I$ (unit matrix), $P(n) = P^n$.

A state is persistent if the probability of sooner or later returning to it is 1, or transient if the probability of sooner or later returning to it is less than 1.

If a Markov chain has a finite number of states, each accessible (probability of going from src to dst is positive) from every other state, then the states are all persistent.

A **continuous Markov process** has similar four conditions to the discrete one introduced before, just that transition probability is defined on continuous time t and don't depend on s .

$$p_{ij}(t) = P(\xi(s+t) = \epsilon_j | \xi(s) = \epsilon_i) \quad i, j = 1, 2, \dots$$

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 $P_{\mathbf{X} \times \mathbf{Y}}$.

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

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

Approximation error is the risk difference between the constrained risk minimizer (in \mathbf{F}) and the target

$$R(f_{\mathbf{F}}) - R(f^*)$$

Estimation error is the risk difference between the constrained empirical risk minimizer and the constrained risk minimizer (both in \mathbf{F})

$$R(\hat{f}_n) - R(f_{\mathbf{F}})$$

Optimization error is the risk difference between a function \tilde{f}_n (which we find in practice) and the constrained empirical risk minimizer *

$$R(\tilde{f}_n) - R(\hat{f}_n)$$

Excess risk is the risk between a function and the target

$$\text{Excess Risk}(\hat{f}_n) = \text{Estimation Error} + \text{Approximation Error} = R(\hat{f}_n) - R(f^*)$$

$$\begin{aligned} \text{Excess Risk}(\tilde{f}_n) &= \text{Optimization Error} + \text{Estimation Error} + \text{Approximation Error} \\ &= R(\tilde{f}_n) - R(f^*) \end{aligned}$$

3.3 Complexity measure, l_1 , l_2 regularization, ridge regression, lasso regression, coordinate gradient descent

Regularization is a mechanism to reduce overfitting. By narrowing down the hypothesis space to only those functions who satisfies being within a complexity measure (or penalize the empirical risk minimizer with another term representing complexity).

Complexity measure $\Omega : \mathbf{F} \rightarrow [0, \infty)$ for linear decision functions $f(x) = w^T x$ can be defined using the p -norm (to the p) of w .

- l_0 complexity is then the number of non-zero coefficients
- l_1 **lasso** complexity is $\sum_{i=1}^d |w_i|$
- l_2 **ridge** complexity is $w^T w = \sum_{i=1}^d w_i^2$

To factor complexity measure in constrained empirical risk minimization, we can use **Ivanov** or **Tikhonov** regularization. where

$$\arg \min_{f \in \mathbf{F}} \frac{1}{n} \sum_{i=1}^n l(f(x_i), y_i) \text{ s.t. } \Omega(f) \leq r$$

is Ivanov regularization, and

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

*Optimization can be negative, however, this $\hat{R}(\tilde{f}_n) - \hat{R}(\hat{f}_n)$ can't. In this hypothesis space, this function can fit the overall world better than the constrained empirical risk minimizer, but it can't on the training set, on which the constrained ERM performs best in this space

is Tikhonov regularization. The two can be shown to be equivalent for many performance measure of f , and Ω , using Lagrangian duality theory. Plugging in the lasso / ridge definition of Ω to either form, we get lasso / ridge regression.

For example, ridge regression in Tikhonov form looks like

$$\hat{w} = \arg \min_{w \in \mathbf{R}^d} \frac{1}{n} \sum_{i=1}^n (w^T x_i - y_i)^2 + \lambda \|w\|_2^2$$

Lasso regression typically lead to feature sparsity (more coefficient equal to 0). This can be shown intuitively by comparing the two on the intersection of contour plot (of empirical risk minimizer) and the area constrained by lasso (diamond on a $f(x) = ax + b$ prediction function) / ridge (circle on the same prediction function) regression, as show in Figure 1*.

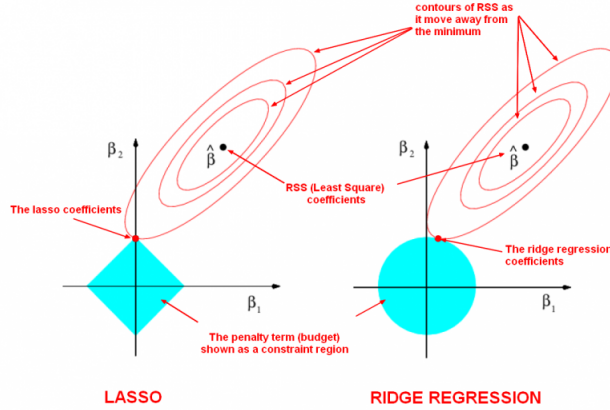


Figure 1: Why lasso regression tends to lead to feature sparsity

Lasso regression (l_1 -norm) is not differentiable, thus we cannot apply gradient descent as-is. We split each coefficient into positive and negative parts: rewrite the vector $w = w^+ - w^-$ ($|w| = w^+ + w^-$), and we have this equivalent problem[†]:

$$\arg \min_{w^+, w^-} \frac{1}{n} \sum_{i=1}^n ((w^+ - w^-)^T x_i - y_i)^2 + \lambda (w^+ + w^-)$$

subject to $w_i^+ \geq 0 \forall i, w_i^- \leq 0 \forall i$

Now we have a constraint and two vectors, to find the minimum, we can use **projected SGD[‡]** or **coordinate descent**.

*<https://niallmartin.wordpress.com/2016/05/12/shrinkage-methods-ridge-and-lasso-regression/>

[†]whose equivalence can be proved. And this can be plugged in to a quadratic solver to give us w^+ and w^-

[‡]Normal gradient descent but project / reset each coordinate to the constrained space after each step

Coordinate gradient descent on a vector w works by adjusting only a single w_i in each step, as opposed to possibly alter the entire w in one step as in regular gradient descent. We iteratively adjust each coordinate several times, where we can pick a random one to adjust, or do cyclic adjustment.

Coordinate gradient descent has a **closed form solution** for lasso (in the form below*), where we can use a formula to solve each w_i , instead of having to loop over them.

$$\hat{w}_j = \arg \min_{w_j \in \mathbf{R}} \frac{1}{n} \sum_{i=1}^n (w^T x_i - y_i)^2 + \lambda \|w\|_1$$

3.4 Elastic nets

Consider what happens in lasso and ridge regressions when features are linear: lasso would assign all the weight to the feature with larger scale, and ridge would assign the weights proportional to the features' scales.

It's similar when features are highly correlated (near-linear relationship between features): think Figure 1 instead of parallel lines, we get elongated ellipses, whose intersection with the hypothesis space would reflect the scales of correlated features.

Elastic net combines lasso and ridge penalties.

$$\hat{w} = \arg \min_{w \in \mathbf{R}^d} \frac{1}{n} \sum_{i=1}^n (w^T x_i - y_i)^2 + \lambda_1 \|w\|_1 + \lambda_2 \|w\|_2^2$$

Geometrically, we end with a hypothesis space like a diamond with edges bulging out (between a circle and a diamond).

3.5 Regression loss functions

A **distance-based loss function** is a loss function ($l(\hat{y}, y) \in \mathbf{R}$) that only depends on the **residual** $r = \hat{y} - y$. Most regression losses are distance-based.

Distance-based losses are translation-invariant: $l(\hat{y} + a, y + a) = l(\hat{y}, y)$.

Square loss ($l(r) = r^2$) penalizes outlier points more heavily than absolute (Laplacian) loss ($l = |r|$) does, and is considered less robust.

Downside with absolute loss is that it's not differentiable.

We are able to construct **Huber loss function** that is robust and differentiable: quadratic for $|r| \leq \delta$ and linear for $|r| > \delta$.

3.6 Classification loss functions

If we have action space \mathbf{A} and outcome space \mathbf{Y} both being $\{-1, 1\}$.

*This form is not differentiable, but it can be shown that for coordinate gradient descent to work there is a weaker condition, which lasso satisfies. *This would indicate that the split into positive and negative shown above is not required, for solving using coordinate gradient descent?*

0-1 loss for $f : \mathbf{X} \rightarrow \{-1, 1\}$:

$$l(f(x), y) = 1 \text{ (} f(x) \neq y \text{)}$$

Where $1 \text{ (} f(x) \neq y \text{)}$ denotes score 1 whenever $(f(x) \neq y)$.

If we allow **real-valued prediction function** $f : \mathbf{X} \rightarrow \mathbf{R}$, meaning action space $\mathbf{A} = \mathbf{R}$ where the value represents confidence of our prediction.

We can define **margin** m as $m = y\hat{y}$, where \hat{y} stands for the predicted score. We want to maximize the margin. Most classification losses are **margin-based**. Empirical risk for 0-1 loss is

$$\hat{R}_n(f) = \frac{1}{n} \sum_{i=1}^n 1 \text{ (} y_i f(x_i) \leq 0 \text{)}$$

$\hat{R}_n(f)$ is non-convex, not differentiable, discontinuous, and its optimization is NP-hard.

SVM/Hinge loss is defined as

$$l_{Hinge} = \max\{1 - m, 0\} = (1 - m)_+$$

It is a convex, upper bound on 0-1 loss, and not differentiable at $m = 1$. And we have **margin error** when $m < 1$.

Using hinge loss function, we define **soft margin linear support vector machine** (a constrained empirical risk minimizer) as

$$\arg \min_{w \in \mathbf{R}^d} \frac{1}{n} \sum_{i=1}^n (1 - y_i f_w(x_i))_+ + \lambda \|w\|_2^2$$

With l_2 regularization, and hypothesis space is linear $\{f(x) = w^T x | w \in \mathbf{R}^d\}$.

Logistic/log loss is defined as

$$l_{Logistic} = \log(1 + e^{-m})$$

It always wants more margin, and is differentiable.

0-1 loss, hinge and Logistic loss functions are illustrated in Figure 2 *.

3.7 SVM and Lagrangian duality

Adding a **bias** term b to the constrained empirical risk minimizer of **support vector machine**, we have

$$\arg \min_{w \in \mathbf{R}^d, b \in \mathbf{R}} \frac{c}{n} \sum_{i=1}^n \max(0, 1 - y_i(w^T x_i + b)) + \frac{1}{2} \|w\|_2^2$$

*<http://fa.bianp.net/blog/2013/loss-functions-for-ordinal-regression/>

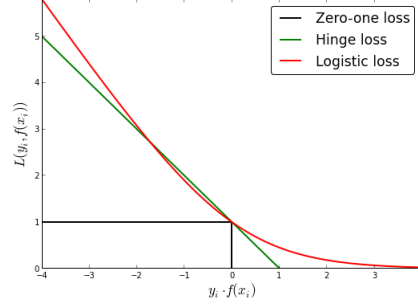


Figure 2: 0-1 loss, hinge loss and Logistic loss functions

whose solution gives us the SVM prediction function. c is a constant that can be tweaked to decide how much regularization matters. This problem is not differentiable due to \max . We can turn it into an equivalent problem.

$$\begin{aligned}
& \text{minimize} && \frac{c}{n} \sum_{i=1}^n \xi_i + \frac{1}{2} \|w\|_2^2 \\
& \text{subject to} && -\xi_i \leq 0, \quad \forall i \in \{1, \dots, n\} \\
& && (1 - y_i(w^T x_i + b)) - \xi_i \leq 0, \quad \forall i \in \{1, \dots, n\}
\end{aligned}$$

This has a differentiable objective function, $n+d+1$ unknowns and $2d$ affine constraints. This can be solved by off-the-shelf QP solver. $\sum_{i=1}^n \xi_i$ now represents the margin loss.

We apply **Lagrangian multiplier** to the constrained optimization problem to get the following

$$\begin{aligned}
L(w, b, \xi, \alpha, \lambda) &= \frac{1}{2} \|w\|_2^2 + \frac{c}{n} \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i (1 - y_i(w^T x_i + b) - \xi_i) - \sum_{i=1}^n \lambda_i \xi_i \\
&= \frac{1}{2} w^T w + \sum_{i=1}^n \xi_i \left(\frac{c}{n} - \alpha_i - \lambda_i \right) + \sum_{i=1}^n \alpha_i (1 - y_i(w^T x_i + b))
\end{aligned}$$

Its **primal** and **dual problems** look like

$$\begin{aligned}
p^* &= \inf_{w, \xi, b} \sup_{\alpha, \lambda \geq 0} L(w, b, \xi, \alpha, \lambda) \\
&\geq \sup_{\alpha, \lambda \geq 0} \inf_{w, \xi, b} L(w, b, \xi, \alpha, \lambda) \\
&= d^*
\end{aligned}$$

This satisfies strong duality by **Slater's constraint qualification**.^{*}
The **Lagrangian dual function** then becomes

$$\begin{aligned} g(\alpha, \lambda) &= \inf_{w, \xi, b} L(w, b, \xi, \alpha, \lambda) \\ &= \inf_{w, \xi, b} \left(\frac{1}{2} w^T w + \sum_{i=1}^n \xi_i \left(\frac{c}{n} - \alpha_i - \lambda_i \right) + \sum_{i=1}^n \alpha_i (1 - y_i (w^T x_i + b)) \right) \end{aligned}$$

This is convex and differentiable: quadratic in w , and linear in ξ_i and b . The global minima is at $\frac{\partial L}{\partial w} = 0$, $\frac{\partial L}{\partial \xi} = 0$, $\frac{\partial L}{\partial b} = 0$.

Solving the three partial derivations (**first order conditions**), we have

$$\begin{aligned} \frac{\partial L}{\partial w} = 0 &\implies w = \sum_{i=1}^n \alpha_i y_i x_i \\ \frac{\partial L}{\partial \xi} = 0 &\implies \alpha_i + \lambda_i = \frac{c}{n} \\ \frac{\partial L}{\partial b} = 0 &\implies \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned}$$

Replacing the w , ξ and b s from the **dual function**, the **dual problem** then becomes

$$\begin{aligned} \sup_{\alpha} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \alpha_i \alpha_j y_i y_j x_j^T x_i \\ \text{s.t.} \quad & \sum_{i=1}^n \alpha_i y_i = 0 \\ & \alpha_i \in \left[0, \frac{c}{n}\right] \quad \forall i \in \{1, \dots, n\} \end{aligned}$$

Note

- Given a solution α^* to dual, primal solution is $w^* = \sum_{i=1}^n \alpha_i^* y_i x_i$.
- w^* is a linear combination of the data x_1, \dots, x_n . The x_i 's corresponding to $\alpha_i^* > 0$ are called **support vectors**.
- $\alpha^* \in \left[0, \frac{c}{n}\right]$, so c controls the maximum weight on each sample. This is **robust**.
- This is a quadratic objective with n unknowns and $n + 1$ constraints.
- Efficient minimization algorithm, sequential minimal optimization, exists.

^{*}Where a convex problem + affine constraints \implies strong duality iff problem is **feasible**. This problem is feasible if we simply let $w = b = 0$ and $\xi_i = 1 \quad \forall i \in \{1, \dots, n\}$

Since we have **strong duality**, **complementary slackness** holds, where the Lagrangian multiplier \times the constraint function is 0.

$$\begin{aligned}\alpha_i^*(1 - y_i f^*(x_i) - \xi_i) &= 0 \\ \lambda_i^* \xi_i &= \left(\frac{c}{n} - \alpha_i^*\right) \xi_i = 0\end{aligned}$$

This means

- $y_i f^*(x_i) > 1 \implies \alpha_i^* = 0, \xi_i^* = 0$
- $y_i f^*(x_i) < 1 \implies \alpha_i^* = \frac{c}{n}, \xi_i^* > 0$
- $y_i f^*(x_i) = 1 \implies \alpha_i^* \in [0, \frac{c}{n}]$
- $\alpha_i^* = 0 \implies \xi_i^* = 0$ (margin loss is 0), so $y_i f^*(x_i) \geq 1$
- $\alpha_i^* = \frac{c}{n} \implies y_i f^*(x_i) \leq 1$
- $\alpha_i^* \in (0, \frac{c}{n}) \implies y_i f^*(x_i) = 1$

Suppose we have an i s.t. $\alpha_i^* \in (0, \frac{c}{n})$, using the complementary slackness conditions we can solve the bias term b^* , whose value stays the same for any choice of i satisfying $\alpha_i^* \in (0, \frac{c}{n})$. *

$$b^* = y_i - x_i^T w^*$$

If there are no such α_i^* s, then we have a degenerate SVM training problem.

3.8 Level sets / contours, sublevel sets and convex optimization

Let $f : \mathbf{R}^d \rightarrow \mathbf{R}$ be a function. A **level set** or **contour line** for the value c is the set of points $x \in \mathbf{R}^d$ for which $f(x) = c$.

A **sublevel set** for the value c is the set of points $x \in \mathbf{R}^d$ for which $f(x) \leq c$.

If $f : \mathbf{R}^d \rightarrow \mathbf{R}$ is **convex**, then sublevel sets are convex. Level sets and superlevel sets of convex functions are not generally convex, hence, the **standard form of convex optimization problem** uses ≤ 0 to express constraints.

The **implicit form of convex optimization problem** goes as

$$\begin{aligned}\text{minimize } & f(x) \\ \text{subject to } & x \in C\end{aligned}$$

where f is a convex function and C is a convex set.

Also, intersection of convex sets is convex.

With numerical error, it's more robust to average over all eligible i and use the mean of resulting b^

3.9 First-order approximation

Suppose $f : \mathbf{R}^d \rightarrow \mathbf{R}$ is differentiable.

We can use **linear** / **first-order approximation** to predict $f(y)$ given $f(x)$ and $\nabla f(x)$

$$f(y) \approx f(x) + \nabla f(x)^T(y - x)$$

Now suppose $f : \mathbf{R}^d \rightarrow \mathbf{R}$ is convex and differentiable.

The linear approximation to f at x is a **global underestimator** of f .

$$\forall x, y \in \mathbf{R}^d, \quad f(y) \geq f(x) + \nabla f(x)^T(y - x)$$

And if $\nabla f(x) = 0$ then x is a global minimizer of f .

3.10 Subgradients

A vector $g \in \mathbf{R}^d$ is a **subgradient** of $f : \mathbf{R}^d \rightarrow \mathbf{R}$ at x if $\forall z$,

$$f(z) \geq f(x) + g^T(z - x)$$

f is **subdifferentiable** at x if \exists at least one subgradient at x .

The set of all subgradients at x is called the **subdifferential**: $\partial f(x)$

f is convex and differentiable $\implies \partial f(x) = \{\nabla f(x)\}$. At any point x , there can be 0, 1, or infinite many subgradients. $\partial f(x) = \emptyset \implies f$ is not convex.

If $0 \in \partial f(x)$, then x is a **global minimizer** of f .

As a reminder, for function $f : \mathbf{R}^d \rightarrow \mathbf{R}$,

- graph of function lives in \mathbf{R}^{d+1}
- gradient, subgradient of f live in \mathbf{R}^d , and
- contours, level sets, sublevel sets are in \mathbf{R}^d

Gradient at x is orthogonal to level set at x . (assuming f continuously differentiable.)

Now to figure out the **direction** on which to do a subgradient descent.

Let $f : \mathbf{R}^d \rightarrow \mathbf{R}$ have a subgradient g at x_0 , hyperplane H orthogonal to g at x_0 must **support** the level set $S = \{x \in \mathbf{R}^d | f(x) = f(x_0)\}$, meaning H contains x_0 and all of S lies on one side of H . The proof of which[†] suggests the following: Points on subgradient g side of H have larger f -values than $f(x_0)$, and points on $-g$ side may not have smaller f -values, meaning $-g$ may not be a descent direction.

^{*}Where local information gives global information!

[†]Using the definition of subgradient, and inner product $g^T \cdot (y - x_0) > 0$ if y is on the side g points in

In **subgradient descent**, suppose we repeatedly step in a negative subgradient direction $x = x_0 - tg$ where $t > 0$ is the step size and $g \in \partial f(x_0)$.

We can prove $-g$ gets us closer to minimizer. Meaning, suppose f is convex, let $x = x_0 - tg$ for $g \in \partial f(x_0)$, and let z be any point for which $f(z) < f(x_0)$, then for small enough $t > 0$, $\|x - z\|_2 < \|x_0 - z\|_2$.

Proof:

$$\begin{aligned}\|x - z\|_2^2 &= \|x_0 - tg - z\|_2^2 \\ &= \|x_0 - z\|_2^2 - 2tg^T(x_0 - z) + t^2\|g\|_2^2 \\ &\leq \|x_0 - z\|_2^2 - 2t(f(x_0) - f(z)) + t^2\|g\|_2^2\end{aligned}$$

Consider

$$g(t) = -2t(f(x_0) - f(z)) + t^2\|g\|_2^2$$

It's a convex quadratic facing upwards, has zeros at $t = 0$ and $t = 2(f(x_0) - f(z))/\|g\|_2^2 > 0$, so

$$g(t) < 0 \quad \forall t \in \left(0, \frac{2(f(x_0) - f(z))}{\|g\|_2^2}\right)$$

Thus we have $-g$ gets us closer to a smaller $f(z)$.

3.11 Features extraction

Mapping an input space \mathbf{X} (sound wave, image, DNA sequence) to a vector \mathbf{R}^d is called **feature extraction** or **featurization**.

A **feature template** is a group of features all computed in a similar way (e.g. *last_three_chars_of(x) = ...*)*.

A **one-hot encoding** is a feature template that always has exactly one non-zero value.

Features can be encoded with an array (good for dense features), or a map (good for sparse features).

Some factors to consider when featurizing

- Non-monotonicity. E.g. temperature as a feature for health prediction, health is not monotonic with temperature. We can transform the input $\Phi x = [1, \{temperature(x) - 37\}^2]^\dagger$, but this would require domain knowledge, instead, we could do $\Phi(x) = [1, temperature(x), \{temperature(x)\}^2]$, where having one extra feature, we increase the flexibility and make it easier to use.[‡]

*With regularization, our resulting prediction function won't be too complicated.

[†]Array representation of features, where 1 is a bias term

[‡]Think less, make the computer do more

- Saturation. E.g. find products relevant to user's query, where given a product x , we have a feature map $\Phi(x) = [1, N(x)]$ where $N(x)$ = number of people who bought x . However at some point x should saturate: a product bought 50000 times is not necessarily 10 times more relevant than a product bought 5000 times, as a linear model would suggest. We could do $\log(1 + N(x))$, or $\arctan(x)^*$ to slow down the growth. Or we could do a discretization, like $\Phi(x) = [1(x < 10), 1(10 \leq x < 100), 1(100 \leq x)]^\dagger$.
- Interaction. E.g. predicting health with weight and height, it's not that they individually matter, but rather weight relative to height. You can add a cross term $h(x)w(x)$, making $\Phi(x) = [1, h(x), w(x), h^2(x), x^2(x), h(x)w(x)]$

A **predicate** of the input space is a function $P : \mathbf{X} \rightarrow \{\text{True}, \text{False}\}$.

What about given features $\Phi(x) = [x_1, x_2]$, we want to classify if a point will fall into a circle in the two-dimensional input space? We could add $x_1^2 + x_2^2$ as another feature.

Similarly, output space can be transformed as well.

We essentially grow the (linear) hypothesis space by adding more features.

3.12 Kernel methods

From a high level, goal is to allow access to huge feature spaces without suffering heavy computational cost.

With featurization, we can map input space $\mathbf{X} \subset \mathbf{R}^n$ to \mathbf{R}^{d^\ddagger} , and the hypothesis space is of affine functions on feature space looks like

$$\mathbf{H} = \{x \rightarrow w^T \Phi x + b | w \in \mathbf{R}^d, b \in \mathbf{R}\}$$

To get expressive hypothesis spaces using linear models, we need high-dimensional feature spaces (say, adding in all $x_i^{k_i}$ where $\sum k_i \leq C$). This complicates computation, and may cause overfitting. Overfitting can be handled with regularization, and kernel methods can help with memory and computational costs.

For example, recall a featurized SVM prediction function, which is the solution to the following

$$\arg \min_{w \in \mathbf{R}^d, b \in \mathbf{R}} \left(\frac{1}{2} w^T w + \frac{c}{n} \sum_{i=1}^n (1 - y_i [w^T \Phi(x_i) + b])_+ \right)$$

whose dual is

*A sigmoid function. Note how this transformation does not need to be convex.

[†]Where if not bucketed carefully, buckets with few items will have coefficients 0, due to regularization. Instead, we could have $\Phi(x) = [1(x \geq 5), 1(x \geq 10), 1(x \geq 100)]$, *whereas a small bucket would fall back to the previous feature.*

[‡]And as it would seem, this optimization makes sense when $n < d$

$$\begin{aligned}
& \sup_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \Phi(x_j)^T \Phi(x_i) \\
& s.t. \sum_{i=1}^n \alpha_i y_i = 0 \\
& \alpha_i \in [0, \frac{c}{n}] \quad \forall i \in \{1, \dots, n\}
\end{aligned}$$

Where Φx only shows up as inner products with other x s.

A method is **kernelized** if inputs only appear inside inner products. $\langle \Phi(x), \Phi(x') \rangle$ for $x, x' \in \mathbf{X}$. The **kernel function** corresponding to Φ and inner product $\langle \cdot, \cdot \rangle$ is

$$k(x, x') = \langle \Phi(x), \Phi(x') \rangle$$

For example, consider quadratic feature map for $x = (x_1, \dots, x_d) \in \mathbf{R}^d$

$$\Phi(x) = (x_1, \dots, x_d, x_1^2, \dots, x_d^2, \sqrt{2}x_1x_2, \dots, \sqrt{2}x_{d-1}x_d)^T$$

has dimension $O(d^2)$, but for any $x, x' \in \mathbf{R}^d$

$$k(x, x') = \langle \Phi(x), \Phi(x') \rangle = \langle x, x' \rangle + \langle x, x' \rangle^2$$

Bringing down the computation cost from $O(d^2)$ to $O(d)$.*

Continuing on with the SVM example, with linear kernel $k(x, x') = x^T x'$, we define the **kernel matrix (Gram matrix)** as such:

$$K = (\langle x_i, x_j \rangle)_{i,j} = \begin{pmatrix} \langle x_1, x_1 \rangle & \dots & \langle x_1, x_n \rangle \\ \vdots & \ddots & \vdots \\ \langle x_n, x_1 \rangle & \dots & \langle x_n, x_n \rangle \end{pmatrix}$$

Then for the standard Euclidean inner product $\langle x_i, x_j \rangle = x_i^T x_j$, we have

$$K = X X^T$$

Where $X = (x_1, \dots, x_n)^T$, an $n \times d$ vector, and K is all the $\Phi(x)$ -terms we are interested in.[†]

And the dual problem becomes

$$\begin{aligned}
& \sup_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j K_{ji} \\
& s.t. \sum_{i=1}^n \alpha_i y_i = 0 \\
& \alpha_i \in [0, \frac{c}{n}] \quad \forall i \in \{1, \dots, n\}
\end{aligned}$$

*Often useful to think of the kernel function as a **similarity score**

[†]Recall with ridge regression, we worked with $X^T X$, which is $d \times d$, meaning we are interested in kernel methods when $n \leq d$

This also gives us the flexibility to change kernels to a different K of $n \times n$ dimension, which can, e.g., correspond to a high dimensional feature space. This is **kernel trick**.

And there are other non-linear kernels, our quadratic example earlier and **polynomial kernel** $k(x, x') = (1 + \langle x, x' \rangle)^M$, which corresponds to a feature map with all monomials up to degree M , the computational cost of the kernel does not grow with M , while plugging in the featurized terms directly the computational cost grows rapidly with M .

And **Radial Basis Function / Gaussian kernel**, $\forall x, x' \in \mathbf{R}^d$

$$k(w, x) = e^{-\frac{\|x - x'\|^2}{2\sigma^2}}$$

*

3.13 Representer theorem and its usage in kernelization

Consider SVM objective function in a more general term,

$$\arg \min_{w \in \mathbf{H}} J(w)$$

And let $J(w)$ be

$$J(w) = R(\|w\|) + L(\langle w, \Phi(x_1) \rangle, \dots, \langle w, \Phi(x_n) \rangle)$$

Where $\|\cdot\|$ corresponds to the inner product on space \mathbf{H} . $R : [0, \infty) \rightarrow \mathbf{R}$ is the regularization term and $L : \mathbf{R}^n \rightarrow \mathbf{R}$ is the loss term.[†]

If $J(w)$ has a minimizer, then it has a minimizer of the form $w^* = \sum_{i=1}^n \alpha_i \Phi(x_i)$.[‡] With this, we don't have to search over the feature space w^* is in, \mathbf{R}^d , but rather having computed the n $\Phi(x_i)$ s, we search over that space of n dimensions. This lets you deal with infinite dimensional feature space.

As an example, a kernelized objective function for SVM looks like

$$\arg \min_{\alpha \in \mathbf{R}^n} R(\sqrt{\alpha^T K \alpha}) + L(K \alpha)$$

Where there is no direct access to $\Phi(x_i)$, and all references are via kernel matrix K .

*Note how this acts like a similarity score, and with this kernel, your featurization dimension d may grow, but to compute the kernel is always $n \times n$ where n is the number of data points.

[†]Each inner product represents a (linear) prediction. Ridge regression and SVM are both of this form. Lasso is not, as l_1 norm does not correspond to an inner product.

[‡]This is the same conclusion we reached at the end of SVM: w^* is a linear combination of features. This can be proved using Pythagorean theorem and projection theorem: do a projection from w^* to M , span of the data

3.14 Performance evaluation

Intuition,

- Always spend some time building a linear baseline model. (lasso / ridge / elastic net regression; l1 / l2 regularized logistic regression or svm)
- Prefer simpler models if performance is the same.
- Consider building an oracle model, which is helpful to get an upper bound on achievable performance. E.g. fit your validation data without regularization. This can give estimate of the approximation error of our hypothesis space.

Confusion matrix for binary classification problem, the 4-cell matrix of true / false positive / negatives.

- **True positive:** predicted positive, and it's right
- **True negative:** predicted negative, and it's right
- **False positive / type 1 error:** predicted positive, and it's wrong
- **False negative / type 2 error:** predicted negative, and it's wrong

Criteria based on the confusion matrix,

- **Accuracy** $(TP + TN)/(TP + TN + FP + FN)$
- **Error rate** $(FP + FN)/(TP + TN + FP + FN)$
- Accuracy + Error rate = 1.
- Accuracy by itself is often not enough, consider a no-information classifier (always produced 0, e.g.) on a distribution that's predominantly one-sided (e.g. 0).
- **Precision** $TP/(TP + FP)$: high precision means low false alarm rate.
- **True positive rate / recall / sensitivity** $TP/(TP + FN)$: high recall means you are not missing many positives.
- **False negative rate / miss rate** $FN/(FN + TP)$
- **False positive rate / fall-out / false-alarm rate** $FP/(FP + TN)$
- **True negative rate / specificity** $TN/(FP + TN)$
- **F_1 score.** $F_1 = 2 \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$. F_β . Weigh towards precision or recall

In a classification problem, the real-valued prediction function (score function) does not have to be thresholded at 0*. Thresholding the score function differently causes the performance metrics to shift. We can threshold it at a value to meet our performance goals, e.g. *recall* > 80%. Measurement curves based on these rates.

- **Precision-recall curve:** as you change the threshold, how precision and recall change.
- **Receiver-operation characteristic (ROC) curve:** as you change threshold, how TP rate and FP rate change.
- **AUC / AUC ROC:** area under the ROC curve.

3.15 Probabilistic modeling

Instead of producing a single value or classification (regression or classification), produce a probability distribution of values, with which we can decide how likely a test data point is. This is applicable in e.g. anomaly detection.

The approach usually entails selecting a particular probability distribution, and fit the samples to derive the parameters of that distribution. (E.g. the mean value of Poisson distribution)

Stratification means partitioning our input data into groups, and treat each group separately. E.g. we stratify each symbol traded into time ranges like 9:30A to 10:30A on every weekday, and train a model for that particular time range.

Bucketing / binning is the opposite of stratification, where we combine natural groups of data into a single group, and train a model for that single group. E.g. we can bucket multiple tech stock symbols into one group, and train one model to predict their overall price movement.

Stratification brings out specificity (also more bias), bucketing smoothes bias out but lose specificity.

3.16 Neural nets

Linear prediction functions, like SVM, ridge, lasso generate feature vector $\Phi(x)$ by hand, and learn parameter vector w from data. A neural net adds **hidden nodes** between the $\Phi(x)$ and score. Multi-layer perceptron, e.g. in this two-layer representation (two affine combinations):

$$h_i = \sigma(v_i^T \Phi(x))$$

$$score = w^T h + b$$

* > 0 predict positive, < 0 predict negative

where σ is a nonlinear **activation function**, and we need to learn the vectors v_i (with a bias term), w , and scalar b . We can add hidden nodes, hidden layers, etc, where > 1 hidden layer is a deep network.

Hyperbolic tangent is a common activation function.

$$\sigma(x) = \tanh(x)$$

Rectified linear function (ReLU) is another. ReLU in practice often works better.

$$\sigma(x) = \max(0, x)$$

One way to think about neural nets is that it's learning the non-linear featurization functions we want to create. (From the hidden layer h_i to *score* looks just like learning in a linear context, and neural nets learn $\Phi(x)$ to h_i as well). Our objective function (with \tanh) is then differentiable w.r.t. all parameters (we can use gradient descent), but not convex. In practice gradient descent seems sufficient.

Universal approximation theorem: a neural network with one (possibly huge) hidden layer can uniformly approximate any continuous function on a compact set iff the activation function is not a polynomial. (the bias term is necessary)

3.17 Multinomial logistic regression

4 Exercises

4.1 Deriving gradient affine form

1. Given $f(w) = c^T w$, $\nabla f(w)$?

$$f'(x; u) = \lim_{h \rightarrow 0} \frac{f(w + hu) - f(w)}{h} = \lim_{h \rightarrow 0} \frac{c^T hu}{h} = c^T u$$

This shows

$$\nabla f(x) = c$$

2. Given $f(w) = w^T A w$, $\nabla f(w)$?

$$\begin{aligned} f'(w; u) &= \lim_{h \rightarrow 0} \frac{f(w + hu) - f(w)}{h} \\ &= \lim_{h \rightarrow 0} \frac{(w + hu)^T A (w + hu) - w^T A w}{h} \\ &= \lim_{h \rightarrow 0} \frac{w^T A w + h w^T A u + h u^T A w + h^2 u^T A u - w^T A w}{h} \\ &= u^T A w + w^T A u \\ &= w^T A^T u + w^T A u \end{aligned}$$

This shows

$$\nabla f(x) = (w^T A^T + w^T A)^T = (A + A^T)w$$

3. Given $f(w) = \|Aw - y\|_2^2$, $\nabla f(w)$?

$$f(w) = \|Aw - y\|_2^2 = (Aw - y)^T (Aw - y)$$

$$\begin{aligned} f'(w; u) &= \lim_{h \rightarrow 0} \frac{(A(w + hu) - y)^T (A(w + hu) - y) - (Aw - y)^T (Aw - y)}{h} \\ &= \lim_{h \rightarrow 0} (A(w + hu) - y)^T Au + (Au)^T (A(w + hu) - y) \\ &= (Aw - y)^T Au + (Au)^T (Aw - y) \\ &= 2(Aw - y)^T Au \end{aligned}$$

This shows

$$\nabla f(x) = 2((Aw - y)^T A)^T = 2A^T(Aw - y) = 2A^T Aw - 2A^T y$$

4. Given $f(w) = \|Aw - y\|_2^2 + \lambda \|w\|_2^2$, express $f(w) = \|Bw - z\|_2^2$?

The gradient of the two need to be the same, using previous result, we need to solve

$$\begin{cases} A^T A + \lambda I &= B^T B \\ A^T y &= B^T z \end{cases}$$

Note the equivalence between extending a matrix and addition, let

$$B = \begin{pmatrix} A \\ \sqrt{\lambda} I_{n \times n} \end{pmatrix} \text{ and } z = \begin{pmatrix} y \\ 0_{n \times 1} \end{pmatrix}$$

written in block-matrix form.

4.2 Recap: linear regression with square loss