

Eric Zhao's

Machine Learning 2023 Fall Note

from Prof. Shengjie Wang

Let $P(B) \neq 0$

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

Independence: $A \perp B$ iff $P(A \cap B) = P(A)P(B)$

$$P(A|B) = P(A)$$

Discrete: $E[X] = \sum_{x \in \text{val}(X)} x p(x)$ (PMF)

continuous: $E[X] = \int_a^b x f_X(x) dx$ (PDF)

$$\text{Var}[X] = E[(X - E[X])^2] = E[X^2] - E[X]^2$$

$$\text{Var}[a] = 0$$

$$\text{Var}[ax] = a^2 \text{Var}[x]$$

Gaussian Distribution

$$X \sim N(\mu, \sigma^2)$$

$$f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$E[X] = \sigma^2 + \mu^2$$

Machine Learning 2023 Fall Note
from Prof. Shengjie Wang

Let $P(B) \neq 0$

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

Independence: $A \perp B$ iff $P(A \cap B) = P(A)P(B)$
 $P(A|B) = P(A)$

Discrete: $E[X] = \sum_{x \in \text{val}(X)} x p_{x|X}(x)$
(PMF)

continuous: $E[X] = \int_a^b x f_X(x) dx$
(PDF)

$$\begin{aligned} \text{Var}[X] &= E[(X - E[X])^2] = E[X^2] - E[X]^2 \\ \text{Var}[a] &= 0 \\ \text{Var}[ax] &= a^2 \text{Var}[x] \end{aligned}$$

Gaussian Distribution

$$X \sim N(\mu, \sigma^2)$$

$$f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$E[X^2] = \sigma^2 + \mu^2$$

Eric Zhao's

Machine Learning 2023 Fall Note

from Prof. Shengjie Wang

Let $P(B) \neq 0$

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

Independence: $A \perp B$ iff $P(A \cap B) = P(A)P(B)$

$$P(A|B) = P(A)$$

Discrete: $E[X] = \sum_{x \in \text{val}(X)} x p(x)$ (PMF)

continuous: $E[X] = \int_a^b x f_X(x) dx$ (PDF)

$$\text{Var}[X] = E[(X - E[X])^2] = E[X^2] - E[X]^2$$

$$\text{Var}[a] = 0$$

$$\text{Var}[ax] = a^2 \text{Var}[x]$$

Gaussian Distribution

$$X \sim N(\mu, \sigma^2)$$

$$f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$E[X] = \sigma^2 + \mu^2$$

Likelihood:

$$L(\theta | X) = P(X | \theta)$$

$$\text{e.g. } p(x_i | \mu, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right)$$

Log-likelihood

$$\ell(\theta | X) = \log(L(\theta | X))$$

$$\text{e.g. } \log p(x | \mu, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2 - \frac{N}{2} \log(\sigma^2) - \frac{N}{2} \log(2\pi)$$

$$\text{MLE: } \hat{\mu}_{ML} = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\hat{\sigma}_{ML}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu}_{ML})^2$$

$$\hat{\mathbb{E}}[\hat{\mu}_{ML}] = \mu \quad \hat{\mathbb{E}}[\hat{\sigma}_{ML}^2] = \frac{N-1}{N} \sigma^2$$

Two Random Variables

Bivariate CDF: $F_{XY}(x, y) = P(X \leq x, Y \leq y)$

Bivariate PMF: $p_{XY}(x, y) = P(X=x, Y=y)$

Marginal PMF: $p_X(x) = \sum_y p_{XY}(x, y)$

Bivariate PDF: $f_{XY}(x, y) = \frac{\partial^2 F_{XY}(x, y)}{\partial x \partial y}$

Marginal PDF: $f_X(x) = \int_{-\infty}^{\infty} f_{XY}(x, y) dy$

Bayes' :

$$P(B|A) = \frac{P(A \cap B) P(B)}{P(A)}$$

$$\mathbb{E}[g(x, Y)] = \iint g(x, y) f_{XY}(x, y) dx dy$$

If x, Y i.d., then $f_{XY}(x, y) = f_X(x) f_Y(y)$

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

$$\text{cov}[X, Y] = \mathbb{E}[(x - \mathbb{E}[x])(Y - \mathbb{E}[Y])] = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y)$$

$= \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$ measure how two variables change together

magnitude and direction of the linear relationship.

$$\text{Var}[X+Y] = \text{Var}[X] + \text{Var}[Y] + \underline{2\text{cov}[X, Y]}$$

$\equiv 0$ when i.d.

Pearson's correlation (normalized covariance)

$$\rho_{X,Y} = \text{corr}(X, Y) = \frac{\text{cov}[X, Y]}{\sigma_X \sigma_Y} \in [-1, 1]$$

measure the linear correlation between two variables

Multiple Random Variables.

joint distribution:

$$F_{x_1 \dots x_n}(x_1 \dots x_n) = P(x_1 \leq x_1, x_2 \leq x_2, \dots, x_n \leq x_n)$$

Joint probability density function.

$$f_{x_1, \dots, x_n}(x_1, \dots, x_n) = \frac{\partial^n F_{x_1, \dots, x_n}(x_1, \dots, x_n)}{\partial x_1 \dots \partial x_n}$$

chain rule:

$$\begin{aligned} f(x_1, \dots, x_n) &= f(x_n | x_1, \dots, x_{n-1}) f(x_1, \dots, x_{n-1}) \\ &= f(x_1) \prod_{i=2}^n f(x_i | x_1, \dots, x_{i-1}) \end{aligned}$$

If x_1, x_2, \dots, x_n are independent.

$$f(x_1, \dots, x_n) = f(x_1) f(x_2) \dots f(x_n)$$

Covariance Matrix

$$\Sigma = \begin{pmatrix} \text{Var}[x_1] & \dots & \text{cov}[x_1, x_n] \\ \dots & \ddots & \dots \\ \text{cov}[x_n, x_1] & \dots & \text{Var}[x_n] \end{pmatrix} = E[(x - E[x])(x - E[x])^T]$$

Multivariate Gaussian Distribution.

$$f(x, \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

 $E[x|Y=y]$ is a function of y $E[x|Y]$ is a random variable.

Law of total expectation:

$$E[E[x|Y]] = E[x]$$

		Mean	Var
Bernoulli (p)	$\begin{cases} p & x=1 \\ 1-p & x=0 \end{cases}$	p	$p(1-p)$
Binomial (n, p)	$\binom{n}{k} p^k (1-p)^{n-k}$	np	$np(1-p)$
Geometric (p)	$\frac{p(1-p)^{k-1}}{k!}$	$\frac{1}{p}$	$\frac{1-p}{p^2}$
Poisson (λ)	$\frac{e^{-\lambda} \lambda^k}{k!}$	λ	λ^2
Uniform (a, b)	$\frac{1}{b-a}$	$\frac{a+b}{2}$	$\frac{(b-a)^2}{12}$
Gaussian (μ, σ^2)	$\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$	μ	σ^2
Exponential (λ)	$\lambda e^{-\lambda x}$	$\frac{1}{\lambda}$	$\frac{1}{\lambda^2}$

K-means

$$\text{objective: } \min_{\{s_1, \dots, s_k\}} \sum_{j=1}^k \sum_{x \in s_j} \|x - c_j\|^2$$

$$\text{without SGD: } c_j = \frac{1}{|s_j|} \sum_{x \in s_j} x$$

$$\text{with SGD: } c_j = (1-\eta) c_j + \eta \underbrace{\sum_{i=1}^n (x_i - c_j)}_{\frac{\partial f(c_1, \dots, c_k)}{\partial c_j}} \quad \text{if } j \in \arg \min_i \|x_i - c_i\|^2$$

Initialization methods:

1) Random k-partition

2) Forgy seed method : Random start from sample points

3) K-means++ i) random start for c_1 ii) for $2 \dots k$, choose c_j at $P = \frac{d_{j-1}(x)}{\sum_{i=1}^n d_{j-1}(x_i)^2}$
next mean is likely to be far from othersMake K-Means Faster : mini-batch due to Parallelization

Limitation of K-Means

- No global optima guaranteed
- behave bad when clusters have varying sizes, different densities or non-spherical shapes

K-Means Analysis: Connection to Gaussian Mixtures

1. Consider a uniform mixture of k spherical Gaussian distributions with means set to centroids of k-means cov matrix is diagonal and $\Sigma = \sigma^2 I$

$$P(x) \propto \frac{1}{k} \sum_{j=1:k} \exp(-\|x - c_j\|^2)$$

$$\text{Proof: } P(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$P(x|\mu, \Sigma) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \cdot e^{-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)}$$

for Gaussian components

we assume equal variance here?

$$\because \text{spherical GD: } \Sigma = \sigma^2 I \quad \therefore |\Sigma| = \sigma^{2d} \quad \Sigma^{-1} = \frac{1}{\sigma^2} I$$

$$P(x) = \frac{1}{k} \sum_{j=1}^k P(x|\mu_j, \Sigma_j)$$

$$= \frac{1}{k} \sum_{j=1}^k \frac{1}{\sqrt{2\pi} \cdot \sigma} (\sigma^2 I)^{\frac{1}{2}} \cdot e^{-\frac{1}{2} (x-\mu_j)^T (\sigma^2 I)^{-1} (x-\mu_j)}$$

$$= \frac{1}{k} \sum_{j=1}^k \frac{1}{\sqrt{2\pi} \cdot \sigma} e^{-\frac{1}{2\sigma^2} \cdot \|x - \mu_j\|^2} \propto e^{\|x - \mu_j\|^2}$$

2.. Let $z(x)$ f $\arg\min_{c_j} \|x - c_j\|^2$ (the centroid closest to x) we have

$$\sum_{j=1}^k \exp(-\|x - c_j\|^2) \approx \exp(-\|x - z(x)\|^2) \text{ due to exponential decay}$$

The contribution of further centroids is negligible

$$3. \therefore P(x) \approx \text{const} \times \exp\left(-\sum_{i=1}^n \|x - z(x(i))\|^2\right)$$

$$\frac{dN}{dt} = -\lambda N$$

Minimizing intra-cluster distances can be viewed as approximations to maximizing the likelihood

Campus Is the inductive bias here? We want the samples for k-means look like Gaussian uniform mixture?

K-nearest neighbor classification

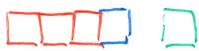
KNN - return the most frequent occurring class of K-nearest neighbors

$$c = \operatorname{argmax}_c \#\{i \mid y_i = c \text{ } i \in \text{NN}(x, k)\}$$

Learning in KNN regression is more about memorizing dataset rather than optimizing model parameters

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - f(x_i))^2$$

Important Rule:



Train V Test

Bayes Classifier:

Assume we know $P(y|x)$, we can predict label by $y^* = \operatorname{argmax} P(y|x)$

Optimal: lowest misclassification rate

No classifier can obtain a lower Bayes error

Given data X and Y , $f(x|Y=c) \sim P_c$ for $c \in \{0, 1\}$

Bayes error $R(h) = \mathbb{E}_X [I_{h(x) \neq Y}]$

$$h^{(B)}(x) = \operatorname{argmax}_{c \in \{0, 1\}} P(Y=c | x=x)$$

$$\text{excess risk: } R(h) - R(h^{(B)}) = 2 \mathbb{E}_X [I_{Z(x) = 0.5} | h_{\text{max}} \neq h^{(B)}(x)]$$

KNN regressor

$$f(x) = \frac{1}{k} \sum_{i \in \text{NN}(x, k)} y_i$$

Make KNN smooth - Kernel

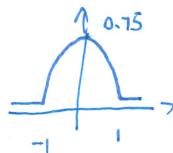
$$f(x) = \frac{\sum_{i=1}^N K_\lambda(x, x_i) y_i}{\sum_{i=1}^N K_\lambda(x, x_i)}$$

Epanechnikov Kernel.

$$K_\lambda(x, x_i) = D\left(\frac{|x-x_i|}{\lambda}\right)$$

$$D(t) = \begin{cases} \frac{3}{4}(1-t^2) & \text{if } |t| \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

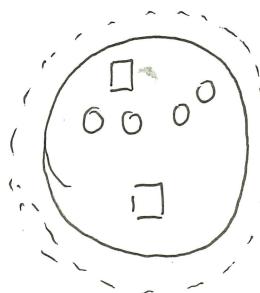
λ determine the width of local neighbourhood



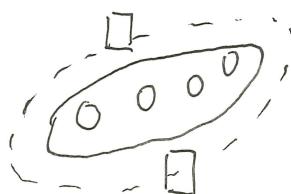
other distance metrics

(弓尺)

Euclidean



Mahalanobis Metric

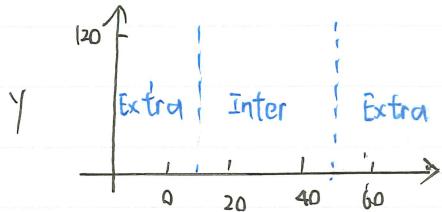


$$\text{Dist}_m(x_1, x_2, \Sigma^{-1}) =$$

$$\sqrt{(x_1 - x_2)^T \Sigma^{-1} (x_1 - x_2)}$$

predict at area without samples

Extrapolation and Interpolation



3. When $d \uparrow$, Euclidean distance is dominated by the biggest difference

$$\lim_{n \rightarrow \infty} E \left(\frac{\text{dist}_{\max}(d) - \text{dist}_{\min}(d)}{\text{dist}_{\min}(d)} \right) = 0$$

4. Samples with Gaussian distribution is most likely to be found at a certain distance from the mean \rightarrow concentrated on the shell

Convex Optimization

convex set C : $\forall x, y \in C, \theta \in [0, 1]$, we have

$$\theta x + (1-\theta)y \in C$$



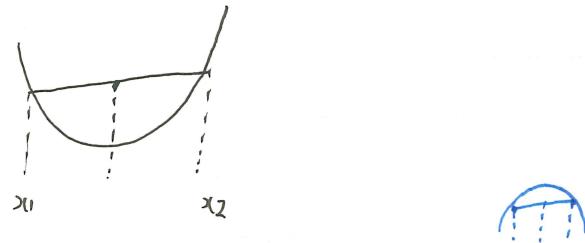
not convex

- \mathbb{R}^d is convex for $\forall d$
- Affine subspace. $\{x | Ax = b\}$. where $A \in \mathbb{R}^{m \times n}$, $b \in \mathbb{R}^m$
- Polyhedron $\{x | Ax \leq b\}$ 多面体
- Intersection of C are C

Convex Functions

$$\forall x_1, x_2 \in D(f) \text{ (convex set)}$$

$$f(\lambda x_1 + (1-\lambda)x_2) \leq \lambda f(x_1) + (1-\lambda)f(x_2) \quad \lambda \in [0,1]$$



- Function f is concave if $-f$ is convex
- convexity = Convexity along all lines
- critical point of convex function is global minimum
- $f'(x)=0$ / $f'(x)$ doesn't exist

3 ways to check convexity

① 0th prove by definition

$$f(\lambda x_1 + (1-\lambda)x_2) \leq \lambda f(x_1) + (1-\lambda)f(x_2)$$

② 1st: tangent line must lie below function.

$$f(x_1) + \frac{d}{dx} f(x_1)(x_2 - x_1) \leq f(x_2)$$

③ 2nd: curvature is always positive

$$\frac{d^2}{dx^2} f(x) \geq 0$$

$$\Delta r_{f(p)} = \lim_{h \rightarrow 0} \frac{f(p+hr) - f(p)}{h}$$

① → ② directional derivative (how function changes in given direction)

① → ③ by calculation and def of derivative

② → ③ expand from $d=0$ (induction)

③ → ④ by mean-value theorem.

Theorem (Stone, 1977)

Let $k \rightarrow \infty$ with increasing n s.t. $\frac{k}{n} \rightarrow 0$ as $n \rightarrow \infty$.

Then for all probability distribution $P(x, y)$ we have $R(\hat{h}_{xy}^{KNN}) \rightarrow R(h^{CB})$

The k should increase to ensure local approximation optimal

To find KNN

- Naive search $O(knd)$ for small d for big d
- (先分区) • tree-based data structures like Kd-tree, ball-tree, LSH
- Approximate nearest neighbour search: locality sensitive hashing
 - similar points have higher chances to be mapped into same bucket

Pros of KNN

- simple, non-parametric
- work well in low dimensions for complex decision surfaces.

Cons of KNN

- computationally expensive, scales with data size.
- suffer from curse of dimensionality

Hyperparameters: k , distance metric, kernel

$$\star \Gamma(n) = (n-1)!$$

Curse of Dimensionality

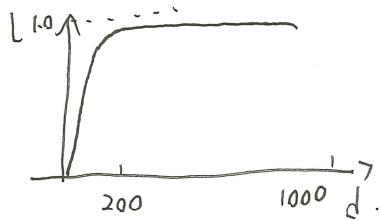
Cause:

- As $d \uparrow$, volume goes higher exponentially
- The space is much more sparse

Phenomena:

1. "Nearest neighbors" could be far away
2. More extrapolation power than interpolation
3. Euclidean distance becomes less meaningful
4. Concentration of Gaussian Distribution lies in peels rather than core

$$1. \quad l = \left(\frac{k}{n}\right)^{\frac{1}{d}} \quad n = \frac{k}{l^d}$$



In 1000d, almost search whole space
for 10NN

2. High dimensional space is more spiky,
All the space is at the edges
The majority of hypersphere tend to concentrate on the shell

$$V_{ball}(d) = V_{ball}(d-1) \cdot \int_{-1}^1 \sqrt{1-x^2}^{d-1} dx$$

$$= \frac{r^d \pi^{\frac{d}{2}}}{r c^{\frac{d}{2}+1}}$$

Checking convexity : 2nd order.

Hessian Matrix positive - semi-definite $\nabla^2 f(x) \geq 0$

A doesn't introduce any "negative weighting"
strong enough to flip the sign. ("weighted dot product")

- Matrix A is positive semi-definite if $x^T A x \geq 0$ for $\forall x \in \mathbb{R}^n$
- A symmetric matrix is p.s.d. iff all eigenvalues are non-negative

$$\nabla^2 f(x) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & & & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

e.g. $H = \begin{bmatrix} 4 & 0 \\ 0 & 2 \end{bmatrix}$



$H = \begin{bmatrix} 4 & 0 \\ 0 & -2 \end{bmatrix}$



$H = \begin{bmatrix} -4 & 0 \\ 0 & -2 \end{bmatrix}$



$H = \begin{bmatrix} 4 & 0 \\ 0 & 0 \end{bmatrix}$



the kernel of SVM must be psd
covariance matrix is psd

For this.
check whether Hessian is
p.d / p.s.d for x
to avoid saddle

Convex Functions:

Exponential: $f(x) = e^{ax}$ for $\forall a \in \mathbb{R}$

Affine function $f(x) = a^T x + b$ where $x \in \mathbb{R}^n$.

both convex and concave, line in 2D and plane in 3D.

Quadratic Functions: $f(x) = \frac{1}{2}x^T A x + b^T x + c$.

for a symmetric matrix $A \in \mathbb{R}^{n \times n}$ and $b \in \mathbb{R}^n$ and $c \in \mathbb{R}$.

Nonnegative weighted sums of convex functions $f(x) = \sum_{i=1}^k w_i f_i(x)$ ≥ 0

Jean's inequality;

Extend from convex function, we have

$$f\left(\sum_{i=1}^k \lambda_i x_i\right) \leq \sum_{i=1}^k \lambda_i f(x_i) \quad \text{for } \sum_{i=1}^k \lambda_i = 1 \text{ and } \lambda_i \geq 0$$

extend to integrals:

$$f\left(\int p(x) x_i dx\right) \leq \int p(x) f(x) dx \quad \text{for } \int p(x) dx = 1 \text{ and } p(x) \geq 0$$

$$f(E[X]) \leq E[f(x)]$$

Gradient Descent

$$x_{t+1} = x_t - \eta \nabla f(x_t)$$

Linear Regression

$$y = w_0 + w_1 x + \epsilon$$

$$MSE = L(w) = \frac{1}{n} \sum_{i=1}^n e_i^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{w}_0 - \hat{w}_1 x_i)^2$$

$$\text{We want } \frac{\partial}{\partial w} L(w) = 0.$$

$$\therefore \hat{w}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}, \quad \hat{w}_0 = \bar{y} - \hat{w}_1 \bar{x}$$

Multivariable:

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - w_0 - w_1 x_{i,1} - \dots - w_p x_{i,p})^2$$

$$= \frac{1}{n} (y - Xw)^T (y - Xw) \quad \text{Convex}$$

$$\text{Goal (differentiate MSE): } -X^T 2(y - Xw) = 0$$

Multicollinearity: one predictor variable in a multiple regression can be linearly predicted from the others.

- sensitive to small changes in data.

- redundant variable

- with multicollinearity, X is not full-rank.

$X^T X$ can't be inverted.

OLS $(X^T X)^{-1} X^T y$ doesn't exist

least square solution not unique

SGD: More iterations, but cheaper

Lower overall cost. Mini-batch SGD achieve better trade-off

Minimizing MSE = Maximizing likelihood

$$L(w) = \log L(w)$$

$$= \sum_{i=1}^n [-\log(\sqrt{2\pi}\sigma) - \frac{(y_i - w^T x_i)^2}{2\sigma^2}]$$

$$= -N \log(\sqrt{2\pi}\sigma) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - w^T x_i)^2$$

Summary of supervised learning

- Assume data is generated i.i.d from unknown distribution P
 $(x_i, y_i) \sim P(x, y)$

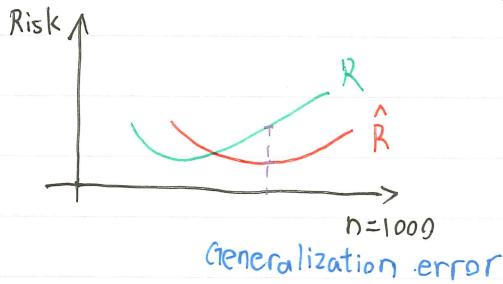
The goal is to minimize expected error (true risk) under P

$$R(w) = \int P(x, y) (y - w^T x)^2 dx dy = E_{x,y} [(y - w^T x)^2]$$

- Estimate the true risk by empirical risk

$$\hat{R}_D(w) = \frac{1}{|D|} \sum_{(x, y) \in D} (y - w^T x)^2$$

We don't want to optimize on training data.



we want to solve $w^* = \arg \min_w R(w)$

not $\hat{w}_D = \arg \min_w \hat{R}_D(w)$

Potential Problem: Outliers, High Leverage Points

$$MSE = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}$$

mean

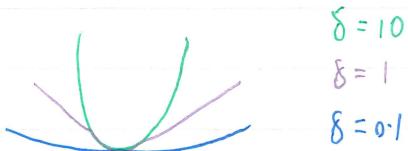
$$MAE = \frac{\sum_{i=1}^n |y_i - \hat{y}_i|}{n}$$

mean

punish less for outliers

Huber Loss (Smooth Mean Absolute Error)

$$L_\delta(y, f(x)) = \begin{cases} \frac{1}{2}(y - f(x))^2 & \text{MSE when } |y - f(x)| \leq \delta \\ \delta|y - f(x)| - \frac{1}{2}\delta^2 & \text{for } |y - f(x)| > \delta \end{cases}$$



Potential Problem - Collinearity

- Consider the level set $\{w | (w - \hat{w})^T X^T X (w - \hat{w}) = c\}$

The equation of an ellipsoid (X full rank), the lengths of the axes scales with eigenvalues $(X^T X)^{-1}$

- When two predictors are highly correlated, the contours of MSE run along a narrow valley, broad range for coefficient estimate
- With x_1 and x_2 linearly related, $X^T X$ has a 0 eigenvalue
- So the level set $\{w | (w - \hat{w})^T X^T X (w - \hat{w}) = c\}$ is no longer an ellipsoid. It's a degenerated ellipsoid, where contour lines are lines pair



hard to get proper coefficients

Decide important variables.

: identify a subset of suitable variables

Forward stepwise selection.

: begins with a model containing no predictors

then add predictors that gives greatest additional improvement

Backward stepwise selection,

• begins with full least squares model,

iteratively removes the least useful predictor

Potential problem: synergy between two predictors

$$y = w_0 + w_1 x_1 + w_2 x_2$$

$$y = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_1 x_2$$

interaction term

Polynomial Regression:

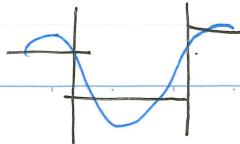
$$y = w_0 + w_1 x + w_2 x^2 + \dots + w_D x^D$$

$$= w^T \phi_1(x) + w_2 \phi_2(x) + \dots + w_D \phi_D(x)$$

$$= w^T \phi(x)$$

Piecewise polynomial:

Represent the function by a separate polynomial in each interval



The function could be discontinuous

$X^T X$: structure of data
 $(X^T X)^{-1}$: structure of model

$$\{b_0\}^2 - 2 \cdot b_0 b_1 \cdot f(x) + f(x)^2 + 2 \cdot f(x) \cdot b_1^2 = f(x) \cdot E(Y)$$

Date _____

Linear Regression:

Compared to the population, our Least Square Estimation

$$\hat{w} = (X^T X)^{-1} X^T Y$$

is a random variable, even though ground truth w^* is not

Bias and Variance Decomposition

suppose the true model is $f^*(x) = E[Y|x]$

Our estimate:

expected prediction over all possible datasets

$$f_D(x) = \arg \min f \sum_{i=1}^n (y_i - f(x_i))^2$$

$$\text{error}(f) = E_p(x) [f(x) - E[Y|x]]^2 - E_p(x, \epsilon) [E[Y|x] - y]^2$$

error in prediction noise, irreducible error, nonrelated

$$E_D[(f_D(x) - E[Y|x])^2] = \underbrace{E_D[(f_D(x) - \bar{f}(x))^2]}_{\text{variance}} + \underbrace{E_D[(\bar{f}(x)) - E[Y|x]]^2}_{\text{bias}}$$

$$E_x(\text{squared loss}) = E_x \{ [E_D(f(x)) - \hat{f}(x)]^2 + \text{Var}_D[\hat{f}(x)] + \text{Var}_\epsilon[\epsilon] \}$$

Assume data is generated using a linear model: $y = X w^* + \epsilon$,

ϵ is random and i.i.d. with mean 0 and variance σ^2

The bias is 0: $E_D[X^T \hat{w}] = X^T w^*$ or $E_D[\hat{w}] = w^*$

The variance is $E_p[(X^T \hat{w} - X^T w^*)^2] = \sigma^2 X^T (X^T X)^{-1} X$

The variance scales with σ^2 , and roughly with cov matrix

$\text{trace}(X^T X)$ \downarrow d and $\uparrow n$ $\uparrow \text{Var} \downarrow \text{Proverfit}$
 \downarrow $d \uparrow \text{Var} \uparrow \text{Pcoverfit}$

roughly eigenvalue of $(X^T X)^{-1}$

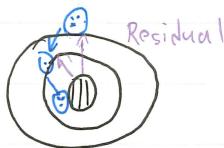
Model complexity \uparrow bias \downarrow variance \uparrow

OLS has smallest error because smallest variance

* bootstrapping : 有放回的多次抽样

Bagging : deal with high variance

Boosting : deal with high bias



in-series model
Boosting



parallel model
Bagging

Ridge Regression:

$$\begin{aligned} \text{In Regression Setting} \\ \min \sum_{i=1}^n (\hat{y}^{(i)} - \sum_{j=1}^p x_j^{(i)} w_j - w_0)^2 \\ \text{s.t. } \sum_{j=1}^p w_j^2 \leq R \end{aligned}$$

$$\begin{aligned} \text{Lagrange Multiplier} \\ \min \sum_{i=1}^n (\hat{y}^{(i)} - \sum_{j=1}^p x_j^{(i)} w_j - w_0)^2 + \lambda \sum_{j=1}^p w_j^2 \end{aligned}$$

- One to one correspondence between λ and R .
- Intercept w_0 is left out of penalty
- Penalization of L2 norm is weight decay

拉格朗日算子:

函数 $f(x, y)$ 与 约束 $g(x, y)$ 相切时有极值，切点法向量平行且梯度方向相同
但梯度大小可能不同，故引入乘子 $\lambda \geq 0$

$$\nabla f(x, y) = \lambda \nabla g(x, y)$$

函数的梯度 = 约束函数的梯度的向量方程

$$\text{Logistic Regression: } \sigma(z) = \frac{1}{1+e^{-z}} = \frac{e^z}{e^z + 1}$$

$$p(x) = \frac{e^{w^T x}}{1+e^{w^T x}}$$

$$\frac{p(x)}{1-p(x)} = e^{w^T x} \quad \text{odds}$$

$$\log\left(\frac{p(x)}{1-p(x)}\right) = w^T x \quad \text{Logit}$$

Maximum Likelihood Estimation.

In MLE, we want to maximize the conditional probability of observing the data given a specific probability distribution and parameters

$$\text{For binary classific.: } P(Y=1|x) = \frac{1}{1+e^{-w^T x}}$$

$$P(Y=0|x) = 1 - \frac{1}{1+e^{-w^T x}}$$

$$\therefore P(Y=y|x) = p(x)^y (1-p(x))^{1-y}$$

$$\max \text{ log-likelihood: } L(w) = \sum_{i=1}^n [y_i \log p(x_i; w) + (1-y_i) \log (1-p(x_i; w))]$$

$$= \min \text{ Cross-Entropy Loss: } L(w) = -\sum_{i=1}^n [y_i \log p(x_i; w) + (1-y_i) \log (1-p(x_i; w))]$$

Entropy: For a discrete event Y : $H(Y=y_0) = -\log(p(y_0))$

For the random variable Y : $H(Y) = -\sum_{i=1}^k p(y_i) \log p(y_i)$

Intuition of quantifying information:

Measuring how much surprise there is in an event

Those events that are rare have more information than those common events

Cross Entropy

$$H(p, q) = - \sum_{i=1}^k p(y_i) \log q(y_i)$$

Expected surprisal of an observer with subjective probabilities Q

upon seeing data that were actually generated according to probabilities P

$H(p, q)$ is minimized when $P = Q$

D_{KL} 散度

$$H(p, q) = H(p) + D_{KL}(p||q)$$

$$D_{KL}(p||q) = - \sum_y p(y) \log \frac{q(y)}{p(y)}$$

$$\frac{\partial L(w)}{\partial w_j} = - \sum_i [y_i - \underbrace{p(x_i; w)}_{\sigma(z)}] x_{ij}$$

Check convexity: $\sigma'(z) = \sigma(z)(1-\sigma(z))$

$$\nabla L(w) = \sum_i x_i x_i^T p(x_i; w) (1 - p(x_i; w))$$

$x_i x_i^T$ is P.S.D. since $(x_i x_i^T)_{ii} = (x_i^T x_i)_{ii} = 1$ and $(x_i^T x_i) \geq 0$

Newton-Raphson algorithm: to find θ s.t $f(\theta) = 0$

$$\theta \leftarrow \theta - \frac{f(\theta)}{f'(\theta)}$$

For optimization at $f(\theta) = 0$

$$\theta \leftarrow \theta - \frac{f'(\theta)}{f''(\theta)}$$

By Taylor expansion: $0 = f(a) = f(\theta^t) + f'(\theta^t)(a - \theta^t) + \frac{1}{2} f''(\theta^t)(a - \theta^t)^2$

$$\frac{f(\theta)}{f'(\theta^t)} + (a - \theta^t) = \frac{-f''(\theta^t)}{2f'(\theta^t)} (a - \theta^t)^2 \rightarrow a - \theta^{t+1} = \frac{-f''(\theta^t)}{2f'(\theta^t)} (a - \theta^t)$$

$$\begin{matrix} a \cdot b \\ c \cdot d \end{matrix}$$

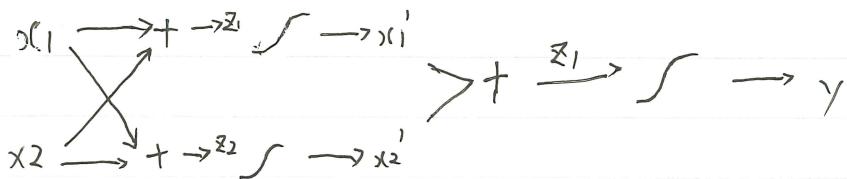
$$\begin{matrix} a \cdot c \\ b \cdot d \end{matrix}$$

$$\begin{matrix} a^2 \cdot b^2 \\ a \cdot c \cdot b \cdot d \end{matrix}$$

No.

Date

Logistic Regression is well calibrated.



Feature Transformation

Classification

$$\begin{aligned} E_D[(f_D(x) - E[Y|x])^2] &= E_D[A^2] + E_D[B^2] \\ &\leq E_D[(\hat{f}_D(x) - \bar{f}(x))^2] + E_D[(\bar{f}(x) - E[Y|x])^2] \end{aligned}$$

Variance

Bias.

Difference : model trained from random data

mean model

mean model

data distribution.

Gauss-Markov Theorem: Among all unbiased estimators, OLS estimator has the smallest variance and hence has the smallest error of all

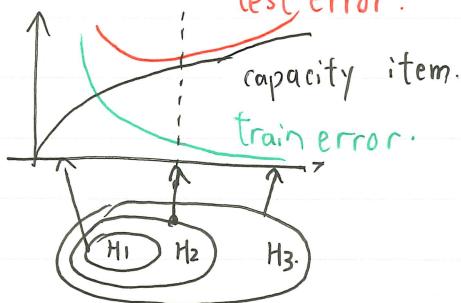
$$\text{Var}(x^T \hat{w}) \leq \text{Var}(\text{Chilnear-unbiased}(x))$$

Complexity Search Space.

test error.

capacity item.

train error.

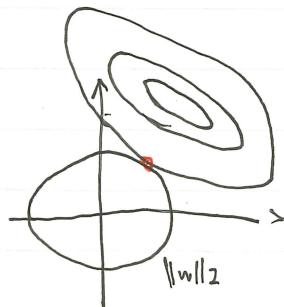


Ridge Regression Lagrange Multiplier

∴ Optimal is achieved when

$$\nabla(L(w) + \lambda \|w\|_2^2) = 0$$

- The point where two contours are tangent to each other

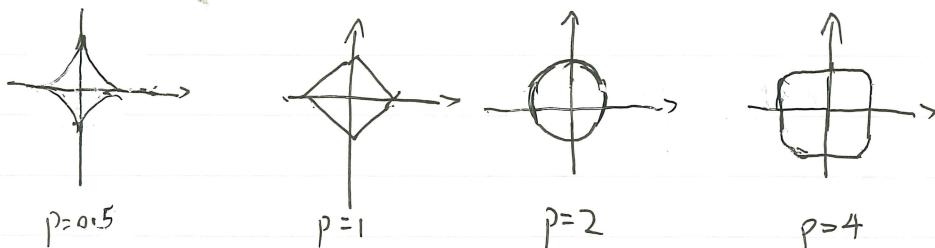


- standard least squares coefficient estimate is **scale-equivalent**
multiply X by $c \rightarrow$ scaling w^T by $\frac{1}{c}$
- Ridge regression is **NOT** scale-equivalent
we usually standardize before apply ridge

close-form solution: $-\frac{1}{n} \sum_{i=1}^n x_i (1_{[y_i=y]} - \frac{e^{x_i}}{e^z}) + \gamma w$

Lasso Regression

$$\min \sum_{j=1}^n (y_i - \sum_{j=1}^d x_{ij} w_j - w_0)^2 + \lambda \sum_{j=1}^d |w_j|$$



$$\|w\|_p = \left(\sum_{j=1}^d |w_j|^p \right)^{\frac{1}{p}}$$

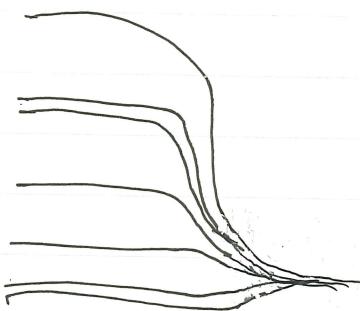
$$\frac{\partial \|w\|_1}{\partial w_j} = \begin{cases} 1 & \text{if } w_j > 0 \\ -1 & \text{if } w_j < 0 \\ ? & \text{if } w_j = 0 \end{cases}$$

Subgradient: For a convex function f , a vector g is a subgradient of f at x . if for all z we have $f(z) \geq f(x) + g^T(z-x)$

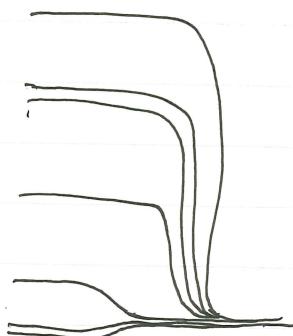
A convex function (not necessarily differentiable) is minimized at x^* . if the zero-vector is a subgradient at x^* .

$$0 \in \partial f(x^*) \text{ iff } x^* \in \operatorname{argmin}_x f(x)$$

Ridge $\|w\|_2$



Lasso $\|w\|_1$



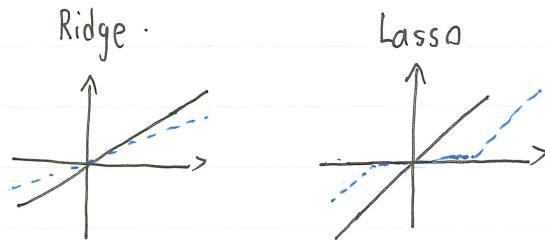
$$w_{ij} = y_i - (\sum_j w_j x_{ij}) , a_j = \sum_i x_{ij}^2, c_j = 2 \sum_i x_{ij} x_{ij}$$

$$w_j = \begin{cases} \frac{c_j + \lambda}{a_j} & \text{if } c_j < -\lambda \\ 0 & \text{if } |c_j| \leq \lambda \\ \frac{c_j - \lambda}{a_j} & \text{if } c_j > \lambda \end{cases}$$

soft thresholding

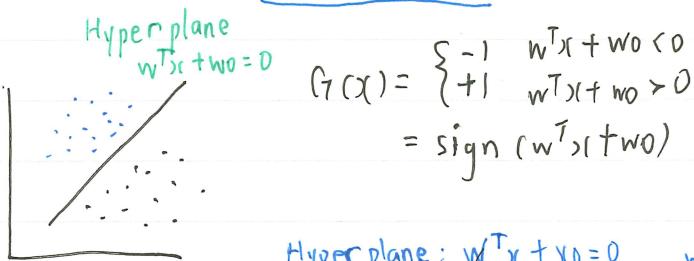
$$\operatorname{sign}(c_j) \max(0, |c_j| - \lambda) / c_j$$

Estimator	Formula
Feature selection	$\hat{w}_j \cdot I(\hat{w}_j > \hat{w}_m)$
Ridge	$\hat{w}_j / (1 + \lambda)$
Lasso	$\text{sign}(\hat{w}_j) (\ \hat{w}_j\ - \lambda) +$ translates coefficient by λ , truncating at 0



SVM (Support Vector Machine)

find the largest margin between classes of points

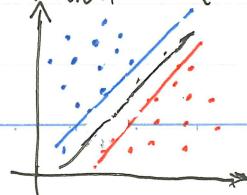


Hyperplane: $w^T x + w_0 = 0$ with normal vector w

Margin: smallest distance from samples to hyperplane

Maximal Margin classifier → maximal margin hyperplane

Support Vectors: 3 training observations that are equidistant from the maximal margin hyperplane
The rest samples are irrelevant.



for Hyperplane $w_0 + w^T x = 0$ normal vector is $\bar{w} = \frac{w}{\|w\|_2}$

Signed distance of any point x to h :

$$\bar{w}^T(x - x_0) = \frac{1}{\|w\|_2} (w^T x + w_0)$$

$$= \frac{1}{\|w\|_2} f(x)$$

Optimization Objective:

$$\max_{w, w_0} M \quad \min_{w, w_0} \text{dist}(x_i, h) \geq m$$

$$\text{s.t. } \forall i: \frac{y_i}{\|w\|_2} (w^T x_i + w_0) \geq M.$$

for one optimal solution w^* and w_0^*

w^* and w_0^* is also an optimal solution.

A Trick: set M to $\frac{1}{\|w\|_2}$

$$\therefore \min \|w\|_2 \quad (\text{maximize } \frac{1}{\|w\|_2})$$

Another trick

$$\min_{w, w_0} \frac{1}{2} \|w\|_2^2$$

$$\text{s.t. } \forall i: y_i (w^T x_i + w_0) \geq 1$$

For Lagrange Multiplier. a more abstract form:

$$\min r(w)$$

$$\text{s.t. } \forall i: g_i(w) \leq 0$$

Lagrangian: $L(w, \alpha) = r(w) + \sum_{i=1}^n \alpha_i g_i(w)$

lagrangian multipliers

g_i and r are convex, and solutions satisfy Karush-Kuhn-Tucker conditions

Primal $\min_{w} \max_{a \geq 0} L(w, a)$

Inner Optimization.

Given a fixed value w . we want to vary a , s.t.
maximize the penalty for violating the constraint

Outer Optimization.

After inner optimization, we vary w to
minimize the loss given the best efforts of a controller

Dual: $\max_{a \geq 0} \min_w L(w, a)$

Karush-Kuhn-Tucker Conditions (KKT)

$$L(w, a) = r(w) + \sum_{i=1}^n a_i g_i(w)$$

$$\left\{ \begin{array}{l} \frac{\partial L}{\partial w}(w, a) = 0 \\ a_i g_i(w) = 0 \\ g_i(w) \leq 0 \\ a_i \geq 0 \end{array} \right. \quad \text{saddle point is optimal}$$

Two scenarios

① at solution $g_i(w^*) < 0$, it's optimized in an unrestricted way,
 $\frac{\partial r}{\partial w}|_{w^*} = 0, a_i^* = 0$

② at $g_i(w^*) = 0, \exists a_i^* \text{ s.t. } \frac{\partial L}{\partial w}(w^*, a_i^*) = \frac{\partial r}{\partial w}|_{w^*} - a_i^* \frac{\partial g_i}{\partial w}|_{w^*} = 0$
 $a_i^* > 0$

Dual Problem of SVM

Primal ~ dual problem

- different perspective
- Easier to solve

$$= \frac{1}{2} \|w\|^2 - \sum_{i=1}^n a_i [y_i (w^T x_i + w_0) - 1]$$

Dual problem: $\max_{\alpha} \min_{w, w_0} L(w, w_0, \alpha)$

Set $\nabla_w L(w, w_0, \alpha) = \nabla_{w_0} L(w, w_0, \alpha) = 0$, we get

$$\begin{cases} w = \sum_{i=1}^n a_i y_i x_i & \textcircled{1} \\ D = \sum_{i=1}^n a_i y_i & \textcircled{2} \end{cases}$$

$$\|w\|^2 = \langle w, w \rangle = \left\langle \sum a_i y_i x_i, \sum a_i y_i x_i \right\rangle \\ = \sum_i a_i y_i \sum_j a_j y_j \langle x_i, x_j \rangle$$

$$= \sum_i \sum_j a_i a_j y_i y_j \langle x_i, x_j \rangle$$

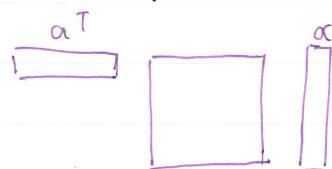
plug in \textcircled{1}, \textcircled{2} $\therefore L(w, w_0, \alpha) =$

$$\sum_i a_i - \frac{1}{2} \sum_i \sum_j a_i a_j y_i y_j \langle x_i, x_j \rangle$$

= A form of Quadratic Programming

$$\min \frac{1}{2} \alpha^T Q \alpha + C^T \alpha$$

s.t. $A \alpha \leq b$



$$Q_{ij} = y_i y_j \langle x_i, x_j \rangle$$

- more general (objective is convex if Q PSD)
- hard to get global minimum for general case.
- Use sequential minimal optimization to solve problem globally
SMO: solve for two coordinates at a time (similar to coordinate descent)

Support Vectors:

From KKT, the constraints $a_i[y_i(w^T x_i + w_0) - 1] = 0$

① $a_i > 0 \quad y_i(w^T x_i + w_0) = 1 \quad x_i \text{ is on the boundary}$

② $y_i(w^T x_i + w_0) > 1 \quad a_i = 0 \quad x_i \text{ is not on boundary}$

$$\begin{aligned} f(x) &= w^T x + w_0 \\ &= w^T \sum_{x' \in [S.V]} y_i a_i \langle x', x \rangle \end{aligned}$$

Line of thought

1. Change constrained optimization into Lagrangian form.
 $L(w, a) = r(w) + \sum_{i=1}^n a_i g_i(w)$ constraints

2. convex \rightarrow primal and dual problem are equivalent

3. study dual problem $\max_{a \geq 0} \min_{w \in W} L(w, a)$

4. KKT conditions \rightarrow simplified to a quadratic programming problem

$$\max \sum_{i=1}^n a_i - \frac{1}{2} \sum_{i,j} y_i y_j a_i a_j \langle x_i, x_j \rangle \quad (\text{to predict based on S.V})$$

5. solve it by quadratic programming solvers

Soft Margin SVM

if two classes are not separable

sample can be on the wrong side of margin or hyperplane

- greater robustness to individual observations
- better classification for most observations

$$\max_{w, w_0, \xi_1, \dots, \xi_n} M$$

subject to $\frac{y_i}{\|w\|_2} (w_0 + w^T x_i) \geq M(1 - \xi_i)$ for $i = 1, \dots, n$.

slack variables $\xi_i \geq 0$ allow observations to lie on the wrong side

$$\sum_{i=1}^n \xi_i \leq \lambda$$

λ hyperparameter

$$\min_{w, w_0} \frac{1}{2} \|w\|_2^2 + C \sum_i \xi_i \quad \text{when } C \uparrow, \text{ more S.V.}$$

subject to $y_i(w_0 + w^T x_i) \geq (1 - \xi_i)$ for $i = 1, \dots, n$

$$\xi_i \geq 0$$

$$\max(1 - y_i(w_0 + w^T x_i), 0) = \xi_i$$

If not violating the margin: $y_i(w_0 + w^T x_i) \geq 1$.

◦ same as the hard margin SVM case. the constant 1 here comes from the trick that we force the norm of w to be $\frac{1}{\lambda}$

◦ Therefore $\xi_i = 0$

If violating the margin, ξ_i indicates how much it violates

if with GD

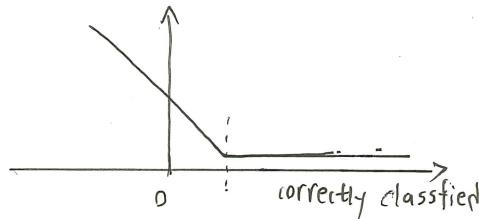
$$\nabla_w L(x_i, y_i) = \left[\mathbb{I}(1 - y_i(w^T x_i + w_0) > 0) y_i \xi_i \right] - w$$

$$\nabla_{w_0} L(x_i, y_i) = \mathbb{I}(1 - y_i(w^T x_i + w_0) > 0) y_i$$

$$\mathbb{I}(Q) = \begin{cases} 1 & Q \text{ is True.} \\ 0 & \end{cases}$$

$$\text{Hinge Loss} \quad \min_{w, w_0} \frac{1}{2} \|w\|^2 + C \sum_i \max(1 - y_i(w_0 + w^T x_i), 0)$$

$$\hookrightarrow L(x_i, y_i, w, w_0) = \begin{cases} 0 & \text{if } 1 - y_i(w_0 + w^T x_i) \leq 0 \\ 1 - y_i(w_0 + w^T x_i) & \text{otherwise} \end{cases}$$



soft Margin SVM uses L2 regularized hinge loss

$$\text{Kernel trick: } \phi(x) = (x_1, x_1^2, \dots, x_1^D, x_2, x_2^2, \dots, x_2^D, x_1 x_2, x_1 x_2^2, \dots, x_1 x_2^{D-1}, \dots)$$

Give $\phi_D(\cdot)$ for polynomial of exactly D-degree.
 $\phi_D(u) \cdot \phi_D(v) = k_D(u, v) = (u \cdot v)^D$

$$\therefore \text{in SVM } f(x) = w_0 + \sum_{(x', y') \in \{s, v\}} y' \langle x', x \rangle$$

$$= \langle w, \phi(x) \rangle + w_0$$

$$= w_0 + \sum_{(x', y') \in \{s, v\}} y' k(x', x)$$

The dual formulation of SVM allow us to use kernel trick \rightarrow capability ↑

Linear kernel: $k(x_i, x_j) = x_i^T x_j$

d-degree polynomial $k(x_i, x_j) = (1 + x_i^T x_j)^d$

Radial: $k(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|^2)$

tiny $k(x_i, x_i)$ when test observation x is far from training observation x_i

Radial basis function kernel is like doing a KNN with smooth neighbourhood

$$\text{let } \gamma = \frac{1}{2} \quad \exp(-\frac{1}{2} \|x_i - x_j\|^2)$$

$$= \exp(-\frac{1}{2} (\|x_i\|_2^2 + \|x_j\|_2^2)) \exp(\langle x_i, x_j \rangle)$$

$$= C \sum_{n=0}^{\infty} \frac{\langle x_i, x_j \rangle^n}{n!}$$

Decision Trees

Tree structure. $T \triangleq (N \cup L, E)$

L : set of leaves N : set of nodes other than other leaves

E : edges of the tree (directed)

$$\Theta = \{T, \{w_j\}_{j \in L}, \{(p_j, T_j)\}_{j \in N}\}$$

$\{w_j\}_{j \in L}$: model used at every leaf

$\{(p_j, T_j)\}_{j \in N}$: splitting function at every node in N

feature dimension to split, splitting threshold

piece-wise constant: constant value in every region.

piece-wise linear: compute $w_j^T x$ in every region

final decision function of decision tree:

$$f(x) = w_q(x)$$

$q(x)$: tree assignment function: it traverse the tree with input x and return the leaf node

when adding 2 more nodes, the objective change from $\min \sum (y_i - w_0)^2$ ($T = \text{mean}(Y)$)
 to $\min_{p_0, T_0} \sum_{(x,y) | x_{p_0} \leq T_0} (y - w_1)^2 + \sum_{(x,y) | x_{p_0} > T_0} (y - w_2)^2$

How to solve?: For every feature dimension p , we sort all data points based on that dimension, and try every middle value as the split, return the best one

When to stop?

1. Max depth / min samples per leaf (hyperparameters)
2. when gain = $\underline{L_{n-1} - L_n}$ is too small

Optimize a tree for classification

Prediction can be a probability vector given the assigned data points

$$w_j = \frac{1}{|R_j|} \sum_{(x,y) \in R_j} \text{onehot}(y, c)$$

a c -length one-hot vector with position $y=1$

Pred class $\leftarrow \arg \max w_{j,i}$

If w_j is one-hot ($\forall i: w_{j,i} \in \{0, 1\}$) all training data points of j has same label,
then the entropy is minimized

If w_j is uniform. $w_{j,i} = \frac{1}{c}$ for all i . the entropy is maximized

Objective for splitting:

$$\max_{P, T} H(w) - \left(\frac{|R_{left}|}{|R|} H(w_{left}) + \frac{|R_{right}|}{|R|} H(w_{right}) \right)$$

Information Gain : $IG(X) = H(Y) - H(Y|X)$

i.e. with more knowledge about the conditioning features.
how much entropy gets reduced

How to solve? Search, find the one with largest decrease

Learning the simplest decision tree is an NP-complete problem

Normally we use greedy approach

1. split current node with largest gain., do this recursively
2. No guarantee on resulting tree.
3. Fast and simple

Good things:

- Easily interpretable
- mirror human decision-making.
- Display graphically
- Trees are **piecewise constant / linear** models.
very capable in fitting the data. (might overfit)

Other tricks: pruning: remove node s.t. decrease in performance is minimal

Random Forest.

- Generate K single trees with randomness, average prediction as final

$$\frac{1}{K} \sum_k f_k(x)$$

Bootstrap: randomly sampling from population to construct new training set

Aggregating (Bagging) Use average output $\frac{1}{K} \sum_k f_k(x)$ as final aggregated predict

Why RF performs better than bagged trees?

- It incorporates a small tweak that decorrelates individual trees
This further reduces variance when we average trees
- The tree may only split on a predictor
from a randomly selected subset of m predictors (features)
- A fresh selection of m randomly selected predictors is presented for each split of each tree
 $m = \sqrt{d}$ turns out to be a good choice.
e.g. we have 100 predictors, each split we choose 10

RF and bagging — ensemble method.

Gradient Boosting Decision Trees

- newly trained tree corrects the error of the previous trees

$$\hat{y}_i^k = \hat{y}_i^{k-1} + f_k(x_i) \quad \begin{array}{l} \text{penalize number of leaves and norm of } w^k \\ \text{regularization} = \gamma |L_k| + \frac{1}{2} \lambda \|w^k\|_2^2 \end{array}$$

$$F_k(f_k) = \sum_{i=1}^n l(y_i, \hat{y}_i^k) + R(\theta_k) = \sum_{i=1}^n l(y_i, \hat{y}_i^{k-1} + f_k(x_i)) + R(f_k)$$

local quadratic approximation

$$\text{with Taylor's expansion. } l(y_i, \hat{y}_i^{k-1} + f_k(x_i)) \approx l(y_i, \hat{y}_i^{k-1}) + g_i f_k(x_i) + \frac{1}{2} h_i f_k^2(x_i)$$

with $k-1$ tree fixed, constant in optimization.

$$F_k(f_k) + \text{const.} \approx \sum_{i=1}^n [g_i f_k(x_i) + \frac{1}{2} h_i f_k^2(x_i)] + \gamma |L_k| + \frac{1}{2} \lambda \|w^k\|_2^2$$

$$= \sum_{i=1}^n [g_i w_{q_k(x_i)}^k + \frac{1}{2} h_i \cdot (w_{q_k(x_i)}^k)^2] + \gamma |L_k| + \frac{1}{2} \lambda \|w^k\|_2^2$$

$w_{q_k(x_i)}^k$: output / weight of leaf node $q_k(x_i)$

w^k is a vector containing all leaf weights.

w_j^k : leaf weight for leaf j

I_j : all the data points assigned to I_j

we change from summing over data point to summing over leaves

$$F_k(\theta_k) + \text{const} \approx \sum_{j \in L_k} \left[C_i \sum_{i \in I_j} q_i w_j^k + \frac{1}{2} (\lambda + \sum_{i \in I_j} h_i) (w_j^k)^2 \right] + \gamma |L_k|,$$

$$\frac{\partial F_k(\theta_k)}{\partial w_j^k} = \left(\sum_{i \in I_j} q_i \right) + \left(\lambda + \sum_{i \in I_j} h_i \right) w_j^k$$

$$\text{for optimal } w_j^k \quad w_j^k = \frac{-G_j}{H_j + \lambda}, \quad G_j \triangleq \sum_{i \in I_j} q_i, \quad H_j \triangleq \sum_{i \in I_j} h_i$$

Now we only need to know I_j , so we need to search over
to set p_j, τ_j

Loss after split on node j :

$$F_k(\theta_k; w_j^k(L), w_j^k(R), \{w_j^k\}_{j \in L_k \setminus j}) = \gamma (|L_k| + 1) - \frac{1}{2} \sum_{j \in L_k \setminus j} \frac{G_j^2}{H_j + \lambda} - \frac{1}{2} \frac{G_j^2(L)}{H_j(L) + \lambda}$$

$$\therefore \text{IG} = \frac{1}{2} \left[\frac{G_j^2(L)}{H_j(L) + \lambda} + \frac{G_j^2(R)}{H_j(R) + \lambda} - \frac{(G_j(L) + G_j(R))^2}{H_j(L) + H_j(R) + \lambda} \right] - \gamma$$

CBDT Algorithm Overview

1. Pick \hat{y}_0 e.g. average of all labels

2. From $k=1$ to K

① For (x_i, y_i) , compute $g_i \triangleq \frac{\partial l(y_i, \hat{y}_i^{k-1})}{\partial \hat{y}_i^{k-1}}$, $h_i \triangleq \frac{\partial^2 l(y_i, \hat{y}_i^{k-1})}{(\partial \hat{y}_i^{k-1})^2} = \frac{\partial g_i}{\partial \hat{y}_i^{k-1}}$

② Build tree k by recursively optimizing improvement current position

by splitting node j

$$\max_{(p_j, \tau_j)} \frac{1}{2} \left[\frac{G_j^2(L)}{H_j(L) + \lambda} + \frac{G_j^2(R)}{H_j(R) + \lambda} - \frac{(G_j(L) + G_j(R))^2}{H_j(L) + H_j(R) + \lambda} \right] - \gamma$$

Boosting: combine multiple weak classifier (generated sequentially)
 use a different loss function or
 reweight training samples (to solve previous mistakes)

Square error: $L(y, f(x)) = (y - f(x))^2 = (e_k y_{i,i} - a g(x_i) y_i)^2$
 learn residual of the current model, more robust to outliers

Exponential loss: $L(y, f(x)) = \exp(-y f(x))$
 $u_i^{(k)} = \exp[-y_i f_{k-1}(x_i)]$ = apply weight to each data point
 make model keep learning after training error is 0

AdaBoost:

Idea: Train $g_2(x)$ on a new training set that fails $g_1(x)$

ϵ_1 : error rate of $g_1(x)$

$$\epsilon_1 = \frac{\sum_{i=1}^n u_i^1 I(y_i \neq g_1(x_i))}{\sum_{i=1}^n u_i^1} \quad \epsilon_1 < 0.5 \quad u_i^k: \text{weight for sample } i \text{ in training classifier } k$$

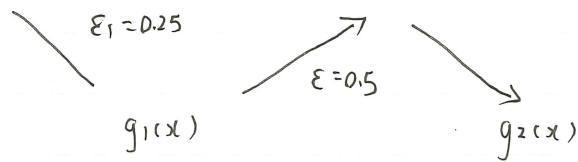
change the example weights from u_i^1 to u_i^2 s.t.

$$\frac{\sum_{i=1}^n u_i^2 I(y_i \neq g_1(x_i))}{\sum_{i=1}^n u_i^2} = 0.5$$

as confusing as possible

Train g_2 based on new weights u_i^2

e.g. $(x_1, y_1, z_1) \quad u_1^{(1)} = 1 \quad \checkmark \quad u_1^{(2)} = \frac{1}{\sqrt{5}}$
 $(x_2, y_2, z_2) \quad u_2^{(2)} = 1 \quad \times \quad u_2^{(2)} = \sqrt{3}$
 $(x_3, y_3, z_3) \quad u_3^{(3)} = 1 \quad \checkmark \quad u_3^{(2)} = \frac{1}{\sqrt{3}}$
 $(x_4, y_4, z_4) \quad u_4^{(4)} = 1 \quad \checkmark \quad u_4^{(2)} = \frac{1}{\sqrt{3}}$



$$g_1(x) = y : \quad u_i^{(2)} \leftarrow u_i^{(1)} / d_1$$

$$g_1(x) \neq y : \quad u_i^{(2)} \leftarrow u_i^{(1)} * d_1 \quad \text{more weight for misclassified classes}$$

let $\bar{\epsilon}_k = \frac{\sum u_i^{(1)}}{d_1} = \sum u_i^{(1)} \cdot d_1$ $d_k = \frac{f_j - \epsilon_k}{\epsilon_k}$
 $g_1(x_i) = y_i \quad g_1(x_i) \neq y_i$

1. Initialize observation weights $u_i^{(1)} = \frac{1}{n}$

2. From $k = 1$ to K .

- Fit classifier g_k to data.
- compute ϵ_k
- update weights

① misclassified $u_i^{(k+1)} \leftarrow u_i^{(k)} \cdot d_k$

② Else: $u_i^{(k)} \leftarrow u_i^{(k-1)} / d_k$

3. Output $f(x) = \text{sign}(\sum_{k=1}^K a_k g_k(x))$

$u_i^{(k)}$

$$\text{for every iteration: } u_i^{(k+1)} = u_i^{(k)} \cdot \exp(-\gamma_i g_k(x_i) \alpha_k)$$

$$\therefore u_i^{(k)} = u_i^{(0)} \prod_{k=1}^{K-1} \exp(-\gamma_i g_k(x_i) \alpha_k)$$

$$= \exp(-\gamma_i \sum_{k=1}^{K-1} g_k(x_i) \alpha_k)$$

$$= \exp(-\gamma_i f_{K-1}(x_i))$$

 α_k

$$\alpha_k = \underset{\alpha}{\operatorname{argmin}} \sum_i e^{-\gamma_i [f_{k-1}(x_i) + \alpha g_k(x_i)]}$$

$L(f_{k-1} + \alpha g_k)$

$$\frac{\partial L}{\partial \alpha} = \sum_i \gamma_i g_k(x_i) e^{-\gamma_i [f_{k-1}(x_i) + \alpha g_k(x_i)]} = 0$$

$$-\sum_{i: g_k(x_i) = 1} e^{-\gamma_i f_{k-1}(x_i) - \alpha g_k(x_i)} + \sum_{i: g_k(x_i) \neq 1} e^{-\gamma_i f_{k-1}(x_i) - \alpha g_k(x_i)} = 0$$

$$-\sum_{i: g_k(x_i) = 1} u_i^{(k)} e^{-\alpha} + \sum_{i: g_k(x_i) \neq 1} u_i^{(k)} e^{\alpha} = 0$$

$$-(1 - \varepsilon_k) e^{-\alpha} + \varepsilon_k e^{\alpha} = 0$$

$$\alpha_k = \frac{1}{2} \ln \frac{1 - \varepsilon_k}{\varepsilon_k} = \ln \left(\sqrt{\frac{1 - \varepsilon_k}{\varepsilon_k}} \right) = \ln(c_k)$$

adaboost factor: $\eta_k = \log \sqrt{\frac{1 - \varepsilon_k}{\varepsilon_k}}$

Naive Bayes (Basically just doing counts)

$$\arg \max_y P(Y=y) \prod_j P(x_j=x_j | Y=y)$$

(If no x_j , add some fake counts)

Point metric:

Label positive	Label negative
----------------	----------------

Predict

Positive

TP

FP

Predict

FN

TN

negative

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

F1-score:

$$F_1 = \frac{2 \cdot \text{Recall} \cdot \text{Precision}}{\text{Recall} + \text{Precision}} \quad (\text{harmonic mean})$$