

Linear Algebra

vector & scalar properties:

- $x+y = y+x$
- $a\vec{0} = \vec{0}$ (zero vector 0)
- $0x = \vec{0}$
- $a(x+y) = ax+ay$
- Distance: $\|x-y\|$

Norms

- Euclidean "2-Norm": $g(x) = \sum_{i=1}^D |x_i|^2$ Manhattan Norm: $g(x) = \sum_{i=1}^D |x_i|$
- Max Norm: $\|x\|_\infty = \max_i |x_i|$
- Zero Norm: $g(x) = \sum_{i=1}^D \mathbb{1}_{\{x_i \neq 0\}}$
- Cauchy-Schwarz Inequality: $|x \cdot y| \leq \|x\| \|y\|$
- Triangle Inequality: adds up non-zeroes

Vector dot product: Orthogonal = perpendicular

$$x \cdot y = \sum_{i=1}^D x_i y_i$$

When $x \cdot y = 0$

[Matrices]

- element-wise product: $\begin{matrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{matrix} \times \begin{matrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{matrix} = \begin{matrix} 16 & 25 & 36 \\ 64 & 80 & 96 \end{matrix}$
- inner prod & dot prod: $x \cdot (y \cdot z) = x \cdot y + x \cdot z$
- (AB)_{n,m} = $\sum_{k=1}^K A_{n,k} B_{k,m}$

Differential Calculus: derivative is the slope of function at that point. often noted to find extreme points, take derivative $\frac{\partial f}{\partial x}$ set to 0. $\frac{\partial f}{\partial x} = 2x^2 - 3x - 1 = 2x(x-3) + 4x - 3 = 0 \Rightarrow x = \frac{3}{4}$

Rules

- scalar multip: $\frac{\partial}{\partial x}[a f(x)] = a [\frac{\partial}{\partial x} f(x)]$ polynomial: $\frac{\partial}{\partial x}[x^k] = kx^{k-1}$
- function addition: $\frac{\partial}{\partial x}[f(x) + g(x)] = [\frac{\partial}{\partial x} f(x)] + [\frac{\partial}{\partial x} g(x)]$
- function mult: $\frac{\partial}{\partial x}[f(x)g(x)] = f(x)[\frac{\partial}{\partial x} g(x)] + [\frac{\partial}{\partial x} f(x)]g(x)$
- function division: $\frac{\partial}{\partial x}\left[\frac{f(x)}{g(x)}\right] = \frac{[\frac{\partial}{\partial x} f(x)]g(x) - f(x)[\frac{\partial}{\partial x} g(x)]}{g(x)^2}$

exponentiation: $\frac{\partial}{\partial x}[e^x] = e^x$ and $\frac{\partial}{\partial x}[a^x] = \log(a)e^x$

logarithms: $\frac{\partial}{\partial x}[\log x] = \frac{1}{x}$

Integral Calculus: computing area under curve

$\int_a^b dx f(x) = \text{area under curve between } a \text{ & } b$
can "kind of" read integral as $\sum_{x=a}^b f(x)$

precision: of predicted true, a high fraction is true, but some true may be predicted false

recall: high recall do not miss true examples, if predicted & false highly likely to be false
but some predicted true may be false

Probability: $P(\text{heads}) = p$ Bernoulli:
 $P(\text{tails}) = 1-p$ $P(x) = p^x(1-p)^{1-x}$

complement: $P(x) = 1 - P(x')$ $x \in \{0, 1\}$

conditional: $P(x|y) = \frac{P(x,y)}{P(y)}$ $P(x|y,z) = \frac{P(x,y,z)}{P(y,z)}$

Addition / Union (Or) ~~*mutually exclusive~~

$P(X \cup Y) = P(X) + P(Y) - P(X \cap Y)$ $P(X \cup Y) = P(X) + P(Y)$

Independence ($X \perp Y$) if:

$P(X|Y) = P(X)$ $P(X,Y) = P(X)P(Y)$

$P(Y|X) = P(Y)$

Product rule: (AND) aka joint probability

$P(X, Y) = P(x) \cdot P(y|x)$ *dependent

$P(x, y, z) = P(x) \cdot P(y|x) \cdot P(z|x, y)$ $P(x, y) = P(x) \cdot P(y)$ independent:

$= P(x) \left[\frac{P(x|y)}{P(y)} \right] \cdot \left[\frac{P(y|z)}{P(z)} \right]$ mutually exclusive:
 $P(x, y) = 0$

Conditional Independence: $(X \perp Y|Z)$

$P(x|y, z) = P(x|z)$

$P(x, y|z) = P(x|z)P(y|z)$

$P(y|x, z) = P(y|z)$

BAYES THEOREM

prior \downarrow likelihood

$$P(C|x) = P(C) P(x|C)$$

posterior

$P(x)$ \leftarrow evidence

posterior: the probability that outcome C_i occurs given some conditions X

prior: what is known regarding possible values that C_i might take before looking at sample or how many times C_i occurred independently of conditions (X)

Likelihood: conditional probability that an event belonging to class C_i has observation X for bag of words

Evidence: # of times outcome X occurred

$$= \frac{P(C_i) P(x|C_i)}{\sum_{k=1}^K P(x|C_k) P(C_k)}$$

for $K=2$ $P(x) = \frac{P(x|C=1)P(C=1)}{P(x|C=1)P(C=1) + P(x|C=0)P(C=0)}$

for something like
 $P(X=1, x_1=0, x_2=1|Y=0)$
 $P(x) = P(x=1, x_2=0, x_3=1|Y=0)$
 $+ P(x=1, x_2=0, x_3=1|Y=1)$

Losses & Risk: $R(x) = \min_i P(C_i|x)$
 $\hat{x}_k = \text{arg} \min_i P(C_i|x)$
 $R(x|x) = \sum_{k=1}^K P(C_k|x) = \sum_{k=1}^K \min_i P(C_i|x)$
 $\text{choose } \hat{x} = \arg \min_i P(C_i|x) = \min_i P(C_i|x)$
 $\text{choose } \hat{x} = \arg \min_i P(C_i|x) = \min_i P(C_i|x)$
 $\hat{x} = \text{average}$

Discriminant Function: maps input variables to class label, parametric classification

- choose C_i if $g_i(x) = m_i x + g_i(x)$
- $g_i(x) = \begin{cases} -1 & C_i \\ P(C_i|x) & P(C_i|x) < P(C_j|x) \end{cases}$
- $P(C_i|x) = \frac{1}{Z} \exp \left[\frac{-(x - \mu_i)^2}{2\sigma_i^2} \right]$
- $Z = \sum_{j=1}^K \exp \left[\frac{-(x - \mu_j)^2}{2\sigma_j^2} \right]$
- $\text{choose } \{C_i \text{ if } g_i(x) > 0\}$
- $\text{choose } \{C_2 \text{ otherwise}\}$

ML estimates

- $P(C_i) = \frac{\sum r_i^i}{N}$ $m_i = \frac{\sum x^i r_i^i}{\sum r_i^i}$
- $\sigma_i^2 = \frac{1}{N} \sum (x^i - m_i)^2 r_i^i$
- $g(x) = \frac{1}{2} \log 2\pi - \log \sigma_i - \frac{(x - \mu_i)^2}{2\sigma_i^2}$ boundary $g_i(x)$

Parametric) Estimator - algorithm (d(x)) to predict i.e. $\hat{\theta} = \sum_{i=1}^N x_i$ parametric
Estimators: θ as random variable w/ prior $p(\theta)$
 Bayes rule: $P(\theta|x) = P(x|\theta)P(\theta)/P(x)$
 MAP - Maximum a posteriori:
 $\theta_{MAP} = \arg\max_{\theta} P(\theta|x) = \arg\max_{\theta} P(x|\theta)p(\theta)$ *posterior estimate*
 Maximum Likelihood (ML): $\theta_{ML} = \arg\max_{\theta} P(X|\theta)$ *w/ infinite data* *MLE & MAP with same result* *and MODE*
 Bayes: $\theta_{Bayes} = \theta_{MAP}$ for Gaussian *mean predictions*
 $E[\theta|X] = \int \theta p(\theta|x)d\theta$ *finds for new data* *MEAN*

- Calculate likelihood:
 For Bernoulli: $f(x) = \prod_{i=1}^N x_i^{x_i} (1-x_i)^{1-x_i}$
 With sample $X = \{x_1, x_2, x_3, \dots, x_N\}$
 $L(\theta|X) = P(x_1|\theta) \cdot P(x_2|\theta) \cdots P(x_N|\theta)$
 $= \theta^{\sum x_i} (1-\theta)^{N-\sum x_i} = \theta^{\sum x_i} (1-\theta)^{N-\sum x_i}$
- Calculate log likelihood:
 $\log L(\theta) = \sum x_i \log \theta + \sum (1-x_i) \log (1-\theta)$
 $= 2 \log(\theta) - \log(1-\theta)$ *row false derivative of log likelihood?*
- For MLE:
 take derivative \rightarrow set to 0
 $L'(\theta) = 2 \theta - 3 \theta^2 = 0$
 $\theta = 0, \frac{2}{3}$ because 0 not in interval

$$E_x[\phi(x)] = E \left[\phi \left(\sum_{i=1}^N X_i \right) \right] = \mu$$

*basically after integration $E_x[\phi(x)]$ = mean

$$\text{Bias } b_e(\theta) = E_x[\phi(x)] - \theta$$

$$\text{Variance: } E_x[(\phi(x) - E_x[\phi(x)])^2]$$

$$\text{Mean square error: } r(\phi, \theta) = E[(\phi(x) - \theta)^2]$$

$$= \text{bias}^2 + \text{variance}$$

↓ bias ↓ variance: Estimator is $\hat{\theta} = \sum_{i=1}^N x_i$; just $\sum x_i$ (not normalized)
 ↓ bias ↑ variance: just $\sum x_i$ (not normalized)
 ↑ bias ↓ variance $\frac{1}{N} \sum_{i=1}^N x_i + E \cdot p^a$ or constant output
 ↑ bias ↑ variance $\frac{1}{N} \sum_{i=1}^N x_i + p^a + E \cdot 10^9$ or sample from very range
 uncertainty or random

$$\text{Likelihood of } \theta \text{ given sample } X \quad (\theta \text{ assumed random var})$$

$$l(\theta|x) = p(x|\theta) = \prod_i p(x_i|\theta) \quad * \text{conditional independence}$$

$$L(\theta|x) = \log l(\theta|x) = \sum_i \log p(x_i|\theta)$$

Maximum Likelihood estimator (MLE)

$$\theta^* = \arg\max_{\theta} L(\theta|x)$$

(3) Find MAP

- take likelihood \times given prior

$$p(\theta|x) = p(\theta) \cdot l(\theta)$$

- plug in prior values

- take log

- false derivative \rightarrow set to 0

Bayes "frequentist estimator"

$$\text{full dist} = p(x_5 | HTHHT) = p(x_5 | \overline{HTHT})$$

$$= \begin{cases} 75\% & x_5 = H \\ 25\% & x_5 = T \end{cases}$$

$$\text{full: } p(x|\theta) = \int p(x|\theta)p(\theta|x)d\theta$$

example: $x_t \sim N(\mu, \sigma^2)$

$$\theta \sim N(\text{C}, \sigma^2)$$

$$\theta_{MLE} = \bar{x} = \sum x_i / N$$

$$\theta_{MAP} = \theta_{BAYES} =$$

$$E[\theta|x] = \frac{N/\alpha}{N/\alpha + 1/\sigma^2} \bar{x} + \frac{1/\sigma^2}{N/\alpha + 1/\sigma^2} \mu$$

sample $\rightarrow N/\sigma^2 + 1/\sigma^2$
 mean \downarrow
 for large samples \downarrow
 prior \uparrow

Distributions:
 Bernoulli: 2 states $x \in \{0, 1\}$
 $p(x_i|\theta) = \theta^{x_i} (1-\theta)^{1-x_i}$
 $\prod p(x_i|\theta) = \prod \theta^{x_i} (1-\theta)^{1-x_i}$
 $MLE \hat{\theta} = \frac{1}{N} \sum x_i$

Multinomial: $K \geq 2$ states $x \in \{0, 1\}$
 $p(x_1, x_2, \dots, x_K | \theta_1, \theta_2, \dots, \theta_K) = \prod_i \theta_i^{x_i}$
 $\prod p(x_i|\theta_i) = \log \prod_i \theta_i^{x_i}$
 $MLE \hat{\theta}_i = \frac{1}{N} \sum x_i$ *because*
be non-negative
add to one

Gaussian (Normal) Dist

$$p(x|\theta) = N(\mu, \sigma^2) \quad \theta = (\mu, \sigma^2)$$

$$p(x|\theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[\frac{(x-\mu)^2}{2\sigma^2} \right]$$

$$\text{MLE for } \mu \text{ and } \sigma^2$$

$$M = \frac{1}{N} \sum x_i \quad S^2 = \frac{1}{N} \sum (x_i - M)^2$$

Dirichlet - multinomial w/ $K \geq 2$ states
 sample likelihood g_i
 $p(x|g) = \prod_{i=1}^K g_i^{x_i} \prod_{i=1}^K \theta_i^{x_i}$
 $\sum_i g_i = 1$

$$\text{prior } q$$

$$\text{Dirichlet}(q|\alpha) = \frac{1}{\Gamma(\alpha_1 \cdots \alpha_K)} \prod_{i=1}^K q_i^{\alpha_i}$$

$$\text{where } \alpha = [\alpha_1, \dots, \alpha_K]^T \quad \alpha_0 = \sum_i \alpha_i$$

$$\Gamma(\alpha) = \text{gamma} = \int_0^\infty u^{\alpha-1} e^{-u} du$$

$$\text{posterior } p(q|x) \propto \prod_{i=1}^K q_i^{\alpha_i + n_i - 1} = \text{Dirichlet}$$

$$\text{where } n = [n_1, \dots, n_K]^T \quad q = \alpha/n$$

$$E[q_i] = \frac{\alpha_i}{\alpha_0 + \sum_i \alpha_i} \quad \text{mode } q_i = \frac{\alpha_i - 1}{\sum_{j \neq i} \alpha_j} \quad \alpha_i > 1$$

Beta variable is binary $x \in \{0, 1\}$
 becomes Bernoulli

$$p(x|q) = \prod_i q^{x_i} (1-q)^{1-x_i}$$

& dirichlet prior reduces to beta:

$$\text{beta}(\alpha, \beta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} q^{\alpha-1} (1-q)^{\beta-1}$$

*conjugate prior

$$\text{posterior: } p(q|A, N, \alpha, \beta) \propto q^{\alpha-1} (1-q)^{N-\alpha+\beta-1}$$

when $\alpha=\beta=1$ uniform & posterior same shape as likelihood

$$\text{mean} = \frac{\alpha}{\alpha+\beta} \quad \text{Naive Bayes assume feature ind. given class}$$

$$\text{mode} = \frac{\alpha-1}{\alpha+\beta-2} \quad p(x_i=1|y) = \frac{\alpha}{\alpha+\beta} \quad \text{given class}$$

$$\text{params} = (1/x_i - 1)$$