

Machine Learning

Supervised

Unsupervised

Active

Collaborative filtering

Machine Translation

out of many approaches in ML, this course will focus on probabilistic (Bayesian) methods

Textbook - Bishop 2006

Terminology

model

sample

likelihood, max likelihood

prior

posterior

MAP

predictive distribution

Model a coin flip

$$P(\text{Head}) = p$$

$$P(\text{Tail}) = 1 - p$$

assumption: no. of flips $\rightarrow \infty$
 \Rightarrow consecutive flips are independent

Model explains how data is generated.

Sample (Sample Data)

H T H H T H T T

x_1, x_2, x_3, \dots
 i^{th} coin flip x_i

$$x_i = \begin{cases} 1 & \text{Head} \\ 0 & \text{Tail} \end{cases}$$

bernoulli variable

$$P(\text{Head}) = P(x_i = 1) = p$$

Scenario 1 200 H 300 T

Scenario 2 2 H 3 T

Scenario 3 15 H 0 T

$$\text{likelihood} = P(\text{data} | \text{model})$$

H T H H T

$$p(1-p) p p (1-p)$$
$$p^3 (1-p)$$

if data $x_1, x_2, x_3, \dots, x_n$

$$P(x_i) = p^{x_i} (1-p)^{1-x_i}$$

$$= \begin{cases} p & x_i = 1 \\ 1-p & x_i = 0 \end{cases}$$

likelihood $L = \prod_{i=1}^n P(x_i)$

$$= \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i}$$

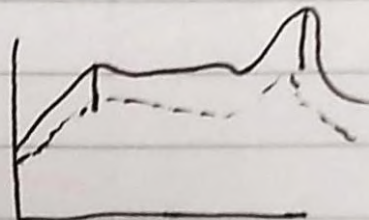
$$= p^{\sum x_i} (1-p)^{n - \sum x_i}$$

$$= p^{\sum x_i} (1-p)^{n - \sum x_i}$$

max likelihood: pick the model that maximizes the value of likelihood.

$$L = p^{\sum x_i} (1-p)^{n - \sum x_i}$$

find p s.t. L is max.



To maximize f

maximize $\log f$ instead (monotonic transformation)

$$L = p^{\#H} (1-p)^{\#T}$$

Take log on both sides

$$\log L = \#H \cdot \log p + \#T \log(1-p)$$

Differentiate wrt p .

$$\frac{d \log L}{dp} = \frac{\#H}{p} + \frac{\#T}{1-p} (-1)$$

set derivative to 0.

$$\frac{\#H}{p} = \frac{\#T}{1-p}$$

$$\frac{\#H}{p} = \frac{\#T}{1-p}$$

$$\#H(1-p) = p \cdot \#T$$

$$\#H - \#H p = p \cdot \#T$$

$$\#H = p(\#H + \#T)$$

$$p = \frac{\#H}{\#H + \#T} = \frac{\#H}{N}$$

max. likelihood estimate of p

$$\hat{p} = \frac{\#H}{N}$$

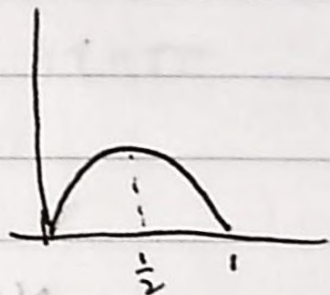
prior = distribution over possible models.

distribution in game over p. choices

choice of prior should reflect our belief/knowledge about the problem.

Choose $p, C(p) \sim C_0 p^2 (1-p)^2$

arbitrarily



$$\int_0^1 C_0 p^2 (1-p)^2 dp = 1$$

$$\left[\frac{p^3}{3} - \frac{(1-p)^3}{3} \right]_0^1$$

$$\int_0^1 p^2 (1+p^2 - 2p) dp = 1$$

$$C \int_0^1 (p^2 + p^4 - 2p^3) dp = 1$$

$$C = 30$$

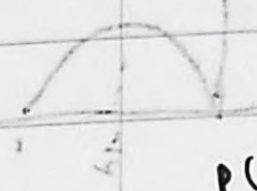
Posterior = distribution over model - that reflects our belief on probability of models, after having seen data.

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

↑ Bayes theorem
↑ model data
↑ P(B)

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

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



$$P(\text{model} | \text{data}) = \frac{\overset{\text{prior}}{P(\text{model})} \overset{\text{likelihood}}{P(\text{data} | \text{model})}}{P(\text{data})}$$

$$\propto P(\text{model}) P(\text{data} | \text{model})$$

$$\propto \text{prior} \times \text{likelihood}$$

here,

$$P(\text{model} | \text{data}) \propto 30p^2(1-p)^2 p^n(1-p)^T$$

$$\propto p^{n+2} (1-p)^{T+2}$$

$$P(\text{Data}) = A \int_0^1 30 p^2 (1-p)^2 p^{200} (1-p)^{200} dp$$

\uparrow All possible ways to generate \uparrow for each test

2nd possible prior (discrete)

P_i	Head		Data
$\frac{1}{2}$	0.9		HTHTT
$\frac{1}{4}$	0.05		
$\frac{3}{4}$	0.05		

$$P(\text{Data}) = 0.9 \binom{5}{2} \left(\frac{1}{2}\right)^5 + 0.05 \binom{5}{1} \left(\frac{1}{4}\right) \left(\frac{3}{4}\right)^4 + 0.05 \binom{5}{4} \left(\frac{3}{4}\right)^4 \left(\frac{1}{4}\right)$$

$p = \frac{1}{2}$ $p = \frac{1}{4}$ $p = \frac{3}{4}$

next flip
Data H T H T H ...

Data H H H H H ...

(next flip)

Prior - distribution over model which expresses our belief (before seeing data) about what are likely models.

choosing arbitrary prior $Pr(P) = 30p^2(1-p)^2$

How to predict without data?

Predicting with data

$$P(\text{model} | \text{Data}) = \frac{P(\text{Data} | \text{model}) P(\text{Model})}{P(\text{Data})}$$
$$\propto P(\text{Data} | \text{model}) P(\text{model})$$

$$\text{Prior} \propto p^2(1-p)^2$$

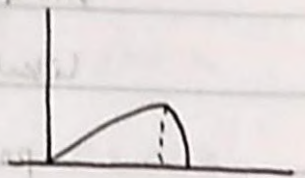
$$\text{Likelihood} = p^3(1-p)^2$$

$$\text{prior} \times \text{likelihood} \propto 15p^5(1-p)^4$$

Situation - H T H T H

NOTE: If we know the functional form of a distribution then we can calculate its normalizing constant.

$$\text{Posterior} = C_2 p^5 (1-p)^4$$



Maximum A Posteriori (MAP): Pick the model (θ) which maximizes the posterior.

$$\text{Posterior} = C_2 p^5 (1-p)^4$$

$$\begin{aligned} \log \text{posterior} &= \log C_2 + \log p^5 + \log (1-p)^4 \\ &= \log C_2 + 5 \log p + 4 \log (1-p) \end{aligned}$$

Differentiate wrt p on both sides

$$\frac{d}{dp} \log \text{posterior} = 0 + \frac{5}{p} + \frac{4}{1-p}$$

set derivative to zero.

$$\frac{5}{p} + \frac{4}{1-p} = 0$$

$$5(1-p) = 4p$$

$$\frac{5}{9} = p$$

vis $\frac{3}{5}$ from max likelihood

situation HHHHH

$$\text{prior } p^2(1-p)^2$$

$$\text{likelihood } = p^5$$

$$\text{posterior} = C_3 p^7(1-p)^2$$

log & d/dp

$$\frac{d}{dp} \log \text{posterior} = \frac{d}{dp} \log C_3 + 7 \frac{d}{dp} \log p + 2 \frac{d}{dp} \log (1-p)$$

$$= 0 + \frac{7}{p} + \frac{2}{1-p}(-1)$$

set derivative to zero.

$$\frac{7}{p} = \frac{2}{1-p}$$

$$7(1-p) = 2p$$

$$\frac{7}{9} = \hat{p}$$

vis $\hat{p} = \frac{7}{9}$ from
maximum
likelihood

Given my current belief: $\Pr(p)$

what is the prob. that next coin is H.

Ex. current belief $\Pr(p) = \begin{cases} 0.9 & p=0.2 \\ 0.1 & p=0.3 \end{cases}$

$$0.9 \times 0.2 + 0.1 \times 0.3 = 0.2$$

Predictive Distribution:

$$P(\text{next head} | \text{data}) =$$

$$\int_P \Pr(p | \text{Data}) P(\text{Next head} | p) dp$$

posterior *likelihood*

Coin Problem:

H Heads

$$\text{Prior} \propto p^2 (1-p)^2$$

T Tails.

$$\text{Posterior} \propto p^{H+2} (1-p)^{T+2}$$

$$\hat{P}_{\text{predictive}} = \int C p^{H+2} (1-p)^{T+2} p dp$$

Probability that $\Pr(\text{next head}) = p$ given data.

prob. that next flip is head given prob head is p .

Conjugate prior : prior and posterior have the same type of dependence on the parameters.

(belong to same family of distributions)

Model - say how data was generated ; prior + dist. over model

Sample - subset of data

method of prediction - max. likelihood

Bayesian posterior

predictive dist.

prior :

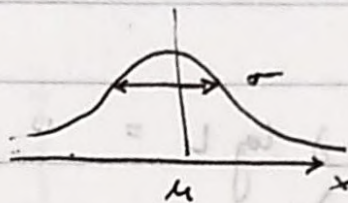
do not rule out options

as the amount of data increases, impact of prior goes down

prior \neq bias

Normal Distribution.

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



precision
 $\beta = \frac{1}{\sigma^2}$

$$f(x) = \frac{\beta}{\sqrt{2\pi}} e^{-\frac{\beta}{2}(x-\mu)^2}$$

Assume that each x_i is drawn from

$N(\mu, \frac{1}{\beta})$ and that different x_i 's are independent

Data: 10.3 5.1 12.6 11.1 2.5

Estimation: pick mean & precision.

$$L = \text{Likelihood} = P(\text{Data} | \text{model})$$

$$= p(x_1) \cdot p(x_2) \cdot p(x_3) \cdot p(x_4) \cdot p(x_5)$$

$$= \prod_i \frac{\beta}{\sqrt{2\pi}} e^{-\frac{\beta}{2}(x_i - \mu)^2}$$

$$\log L = \sum \log \left(\frac{\beta}{\sqrt{2\pi}} e^{-\frac{\beta}{2}(x_i - \mu)^2} \right)$$

$$= \sum \left[-\frac{1}{2} \log 2\pi + \frac{1}{2} \log \beta - \frac{\beta}{2} (x_i - \mu)^2 \right]$$

$$= \frac{N}{2} \log 2\pi + \frac{N}{2} \log \beta - \frac{\beta}{2} \sum (x_i - \mu)^2$$

$$\frac{\partial \log L}{\partial \mu} = 0 + 0 - \beta \sum_i z(x_i - \mu) (-1) = 0$$

$$\sum_i x_i = \sum_i \mu = N\mu$$

$$\hat{\mu}_{ML} = \frac{\sum x_i}{N} \quad \left\| \begin{array}{l} \text{max likelihood} \\ \text{estimator} \end{array} \right.$$

Estimator is a function of random variable(s)

\Rightarrow it is a random variable

- ① unbiased $E[\hat{\mu}] = \mu$
 ② Prefer low variance
- } properties of good estimator

$$\frac{\partial \log L}{\partial \beta} = \frac{N}{2} \frac{1}{\beta} - \frac{1}{2} \sum (x_i - \mu)^2 = 0$$

$$\frac{1}{\beta_{ML}} = \frac{1}{N} \sum (x_i - \mu)^2$$

↑ replace with $\hat{\mu}_{ML}$

$$\frac{1}{\beta_{ML}} = \frac{1}{N} \sum (x_i - \mu)^2$$

↑ biased estimator

$$E(\hat{\mu}_{ML}) = \frac{N\mu}{N} = \mu$$

$$\frac{1}{\beta} = \frac{1}{N-1} \sum (x_i - \hat{\mu}_{ML})^2$$

↑ corrected unbiased

$$E[\hat{\mu}_n] = E\left[\frac{1}{n} \sum x_i\right]$$

$$= \frac{1}{n} \sum E[x_i]$$

$$= \frac{1}{n} \cdot \sum \mu.$$

$$= \frac{1}{n} \cdot n \mu = \mu$$

∴ it is an unbiased estimator

Machine Learning

- ① A model explains how data is generated
- ② We estimate the concrete model (parameters) using data

① Maximum Likelihood - non bayesian

② Prior + data \rightarrow posterior \rightarrow MAP

all information \rightarrow Predictive distribution } bayesian

prior protects against overfitting

coin model : p

$$L = p^{\#H} (1-p)^{\#T}$$

$$\text{prior } \beta(a, b) = \frac{p^{a-1} (1-p)^{b-1}}{\int_0^1 p^{a-1} (1-p)^{b-1} dp} = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$$

Used $\beta(3, 3)$ in example

$$\text{posterior} \propto \text{prior} \times L \propto p^{\#H + a - 1} (1-p)^{\#T + b - 1}$$

$$\text{dist.} \propto p^2 (1-p)^2$$

must be

$$Pr(p) = \frac{\Gamma(3+3)}{\Gamma(3)\Gamma(3)} p^2 (1-p)^2 = \frac{6!}{3!3!} p^2 (1-p)^2$$

$$= \frac{720}{36} p^2 (1-p)^2$$

$$= 20 p^2 (1-p)^2$$

$$\int p^3 (1-p)^7 dp \times \frac{\Gamma(4+8)}{\Gamma(4)\Gamma(8)} = \frac{\Gamma(6)\Gamma(8)}{\Gamma(4+8)}$$

$$\downarrow$$

$$\frac{\Gamma(4)\Gamma(8)}{\Gamma(4+8)} \int \beta(4,8) dp$$

$$\frac{4!7!}{12!}$$

Normal distribution

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

$$= \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}x^2} e^{\frac{x\mu}{\sigma^2}} e^{-\frac{\mu^2}{2\sigma^2}}$$

$$\propto e^{\frac{2\mu x - x^2}{2\sigma^2}}$$

"completing the squares"

$$p(x) \propto e^{-16x^2 + 8x}$$

← normal distribution

$$\frac{2\mu x - x^2}{2\sigma^2} = -16x^2 + 8x$$

$$\frac{+1}{2\sigma^2} = +16 \quad \text{and} \quad \frac{\mu}{\sigma^2} = 8$$

$$8/32 = \frac{1}{4}$$

$$\sigma^2 = 32$$

$$\mu = 8 \times 32$$

Model 2: independent samples from Normal($\mu, \frac{1}{\beta}$)

$$L = \prod_i \sqrt{\frac{\beta}{2\pi}} e^{-\frac{\beta}{2}(x_i - \mu)^2}$$

$$\hat{\mu}_{ML} = \frac{\sum x_i}{N}$$

What would a conjugate prior for μ look like?

prior \times L \propto posterior

$$L = \prod_i \sqrt{\frac{\beta}{2\pi}} e^{-\frac{\beta}{2}(x_i - \mu)^2}$$

$$= \prod_i \sqrt{\frac{\beta}{2\pi}} e^{-\frac{\beta}{2} x_i^2} e^{-\beta \mu \sum x_i} e^{\frac{\beta}{2} \mu^2 \sum 1}$$

exponential quadratic μ

Prior $P(\mu | \mu_0, \beta_0) = \sqrt{\frac{\beta_0}{2\pi}} e^{-\frac{\beta_0}{2}(\mu - \mu_0)^2}$

$$= \sqrt{\frac{\beta_0}{2\pi}} e^{-\frac{\beta_0}{2}\mu^2 + \beta_0 \mu \mu_0 - \frac{\beta_0 \mu_0^2}{2}}$$

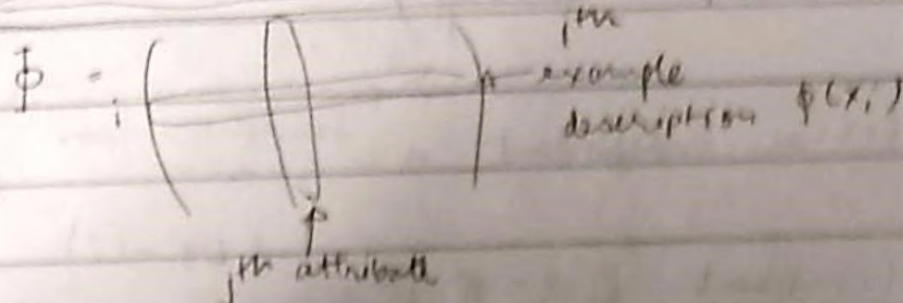
$$= \sqrt{\frac{\beta_0}{2\pi}} e^{-\frac{\beta_0 \mu_0^2}{2}} \left(e^{-\frac{\beta_0}{2}\mu^2 + \beta_0 \mu \mu_0} \right)$$

What are mean & precisions ($= \frac{1}{\sigma^2}$) for posterior?
 (Todo)

House

	Age	Size	Condition	Centrality	Price
House 1	2	1000	10	-1000	10
House 2	100	500	8	1	15
					?

output



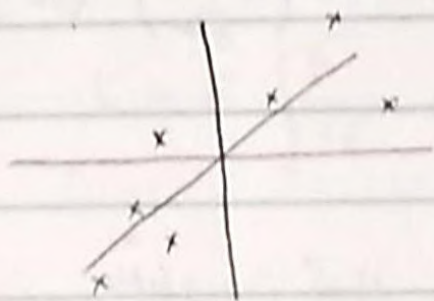
$t =$ vector of values we want to predict,
 one for each example.

NOTE: All vectors are column vectors

$x_i = i^{\text{th}}$ example

$\phi(x_i) =$ representation of i^{th} example

Linear regression



x_i 's are arbitrary

$$y_i = \bar{w}^T \phi(x_i)$$

$$y_i = \sum_j w_j \phi_j(x_i)$$

True

$$t_i \sim N(y_i, \frac{1}{\beta})$$

observed

t_i 's are independent of each other
assume variance of all examples is same.

$$\text{Likelihood } L = \prod_i \frac{\beta}{\sqrt{2\pi}} e^{-\frac{\beta}{2}(y_i - t_i)^2}$$

$$L = \prod_i \frac{\beta}{\sqrt{2\pi}} e^{-\frac{\beta}{2}(w^T \phi(x_i) - t_i)^2}$$

$$\log L = \frac{N}{2} \log \frac{\beta}{2\pi} - \frac{\beta}{2} \sum_{i=1}^N (w^T \phi(x_i) - t_i)^2$$

Max. likelihood for w

$$\max \log L \text{ is same as } \max - \sum_{i=1}^N (w^T \phi(x_i) - t_i)^2$$

$$\text{or } \min \sum_{i=1}^N (w^T \phi(x_i) - t_i)^2$$

A matrix vector product is the linear combination of the column vectors of matrix.

$$\min \sum (w^T \phi(x_i) - t_i)^2 \quad \text{Least Squares.}$$

$$\begin{bmatrix} | & | & | \\ A & & \\ | & | & | \\ a_i & & \end{bmatrix} \begin{bmatrix} | \\ b \\ | \end{bmatrix} = \begin{bmatrix} | & | \\ \sum a_i b_i & \\ | & | \end{bmatrix}$$

$$\begin{pmatrix} a & b & c \\ d & e & f \\ g & h & i \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \\ \gamma \end{pmatrix} = \begin{pmatrix} a\alpha \\ d\alpha \\ h\alpha \end{pmatrix} + \begin{pmatrix} b\beta \\ e\beta \\ i\beta \end{pmatrix} + \begin{pmatrix} c\gamma \\ f\gamma \\ j\gamma \end{pmatrix}$$

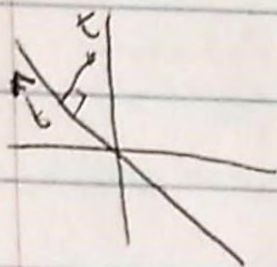
$$\left(\begin{array}{l} \leftarrow \\ \leftarrow \end{array} \right) \begin{array}{l} w^T \phi(x_i) \\ = \phi(x_i)^T w \end{array}$$

$$\begin{aligned} \min \sum (w^T \phi(x_i) - t_i)^2 \\ \min \| \Phi w - t \|^2 \end{aligned}$$

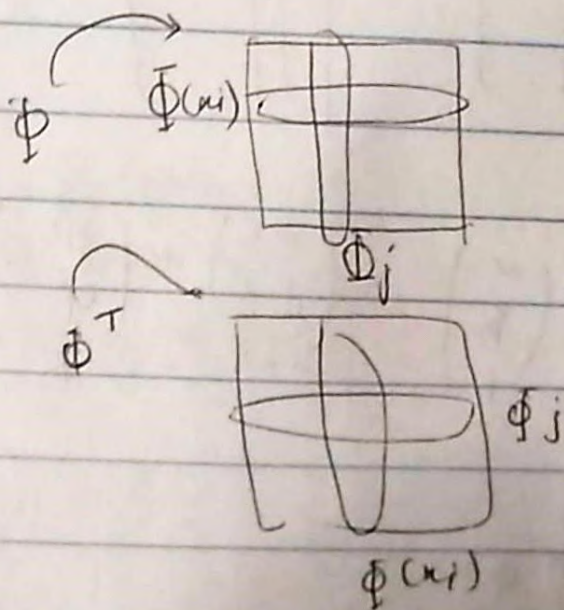
Find \hat{w} such that $\| \hat{t} - t \|^2$ is

minimized where $\hat{t} = \Phi \hat{w}$ prediction parameter

\hat{t} is a linear combination of columns of Φ



~~t~~ + \hat{t} is perpendicular to columns of Φ



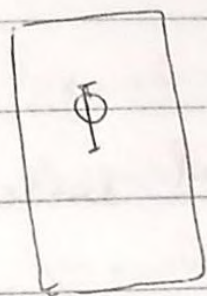
$$\Phi^T (t - \hat{t}) = 0$$

$$\Phi^T (t - \Phi \hat{w}) = 0$$

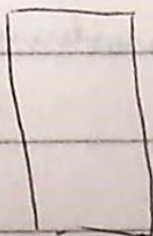
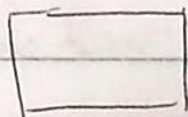
$$\cancel{\Phi^T t} = \Phi^T \Phi \hat{w}$$

$$\Phi^T t = \Phi^T \Phi \hat{w}$$

$$\hat{w} = (\Phi^T \Phi)^{-1} \Phi^T t$$



$$\Phi^T \Phi$$



=



correlation
b/w elements

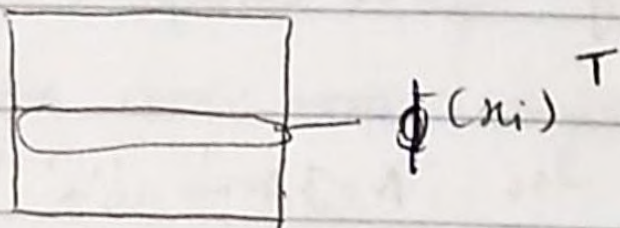
must be invertible.

For any A

$A^T A$ has an inverse

iff. col of A are independent

Machine Learning.



x_i example

$\phi(x_i)$ column vector, represents the example

true hidden value

$$y_i = w^T \phi(x_i)$$

observe

$$t_i \sim N(y_i, \frac{1}{\beta})$$

Every example is drawn independently

$$L = \prod_i \frac{\beta}{\sqrt{2\pi}} e^{-\frac{\beta}{2} (\sum w^T \phi(x_i) - t_i)^2}$$

⋮

Maximum likelihood :

$$\text{Max } \frac{N}{2} \log \beta - \frac{\beta}{2} \sum_{i=1}^N (w^T \phi(x_i) - t_i)^2$$

max L for w

is same as

$$\min_i \sum (w^T \phi(x_i) - t_i)^2$$
$$\min \| \Phi w - t \|^2$$

Soln #1 Geometric

$$\hat{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T t$$

Soln #2 Using Derivatives

$$\min \sum (w^T \phi(x_i) - t_i)^2$$

$$w^T \phi(x_i) = \sum_{k=1}^K w_k \phi_{ik}(x_k)$$

$$\frac{\partial}{\partial w} \sum_i (w^T \phi(x_i) - t_i)^2 = \sum_i 2(w^T \phi(x_i) - t_i) \phi(x_i)$$

= let A

$$\frac{\partial A}{\partial w} = \begin{pmatrix} \frac{\partial A}{\partial w_1} \\ \frac{\partial A}{\partial w_2} \\ \vdots \\ \frac{\partial A}{\partial w_K} \end{pmatrix} = 2 \sum (w^T \phi(x_i) - t_i) \phi(x_i) = \vec{0} = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

what is the vector whose entries are $\phi_k(x_i)$: $\begin{pmatrix} \phi_1(x_i) \\ \vdots \\ \phi_k(x_i) \end{pmatrix}$

$\phi(x_i)$

$$\sum (w^T \phi(x_i)) \phi(x_i) = \sum t_i \phi(x_i)$$

$$\phi^T \phi w \quad \phi^T t$$

vector whose entries are $\sum w^T \phi(x_i) = \sum \phi(x_i)^T w$

59

$$a^T b = \sum_k a_k b_k$$

$$\Rightarrow \Phi^T \Phi w = \Phi^T t$$

$$\Rightarrow \hat{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T t$$

Soln #3 Vector derivatives

$$\text{Min } \|\Phi w - t\|^2$$

$$\|\Phi w - t\|^2 = (\Phi w - t)^T (\Phi w - t) \quad \text{scalar}$$

$$= w^T \Phi^T \Phi w - w^T \Phi^T t - t^T \Phi w + t^T t$$

$$= w^T \Phi^T \Phi w - 2w^T \Phi^T t + t^T t$$

(Assuming A is symmetric)

$$\left\| \frac{\partial w^T A w}{\partial w} = 2Aw \right.$$

$$\left\| \frac{\partial w^T b}{\partial w} = b \right.$$

$$\frac{\partial \|\Phi w - t\|^2}{\partial w} = 2\Phi^T \Phi w - 2\Phi^T t = 0$$

$$\Rightarrow \Phi^T \Phi w = \Phi^T t$$

$$\Rightarrow \hat{w} = (\Phi^T \Phi)^{-1} \Phi^T t$$

Regularization : trick to restrict w from getting arbitrary values & therefore overfitting the data.

New objective : $\min \| \Phi w - t \|^2 + \lambda \| w \|^2$

regularization parameter

penalty for increase in norm of w .

$$\lambda \| w \|^2 = \lambda w^T w$$

$$\frac{d \lambda \| w \|^2}{d w} = 2 \lambda w$$

$$\therefore \frac{d \| \Phi w - t \|^2 + \lambda \| w \|^2}{d w}$$

$$= 2 \Phi^T \Phi w - 2 \Phi^T t + 2 \lambda w$$

$$\frac{d \text{objective}}{d w} = 0$$

$$\Phi^T \Phi w + \lambda w - 2 \Phi^T t = 0$$

$$\Phi^T \Phi w + \lambda I w - 2 \Phi^T t = 0$$

$$(\Phi^T \Phi + \lambda I) w = \Phi^T t$$

$$\hat{w}_R = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T t$$

$$\log L = -\frac{N}{2} \log 2\pi + \frac{N}{2} \log \beta - \frac{\beta}{2} \mathbf{1}^T \Phi \mathbf{W} - \mathbf{t} \mathbf{1}^2$$

compute MLE for β

$$\frac{d \log L}{d \beta} = \frac{N}{2\beta} - \frac{1}{2} \mathbf{1}^T \Phi \mathbf{W} - \mathbf{t} \mathbf{1}^2$$

$$\sigma^2 = \frac{1}{\beta} = \frac{1}{N} \mathbf{1}^T \Phi \mathbf{W} - \mathbf{t} \mathbf{1}^2$$

conjugate prior for \mathbf{W} ?

Prior \times Likelihood \rightarrow posterior

$$L \propto \pi \left(\frac{\beta}{2\pi} \right)^{\frac{1}{2}} e^{-\frac{\beta}{2} (\mathbf{w}^T \Phi(\mathbf{x}_i) - t)^2}$$

$$\propto e^{-\frac{\beta}{2} \mathbf{1}^T \Phi \mathbf{W} - \mathbf{t} \mathbf{1}^2}$$

$$\propto e^{-\frac{\beta}{2} \mathbf{1}^T \Phi \mathbf{W} - \mathbf{t} \mathbf{1}^2} e^{\frac{\beta}{2} (\mathbf{w}^T \Phi^T \mathbf{W} - 2 \mathbf{w}^T \Phi^T \mathbf{t} + \mathbf{t}^T \mathbf{t})}$$

$$\propto e^{-\frac{\beta}{2} \mathbf{w}^T \Phi^T \mathbf{W}} e^{\frac{\beta}{2} \mathbf{t}^T \Phi^T \mathbf{t}} e^{-\frac{\beta}{2} \mathbf{t}^T \mathbf{t}}$$

$$\propto e^{-\frac{\beta}{2} \mathbf{w}^T \Phi^T \mathbf{W}} e^{\frac{\beta}{2} \mathbf{t}^T \Phi^T \mathbf{t}} \quad \text{ignore}$$

$e^{-\frac{\beta}{2} \mathbf{w}^T \Phi^T \mathbf{W}}$ is quadratic
 $e^{\frac{\beta}{2} \mathbf{t}^T \Phi^T \mathbf{t}}$ is linear

~~Note~~ Multivariate Gaussian Distribution

↳ looks like a gaussian in n -dimensions

* Read Linear Algebra review

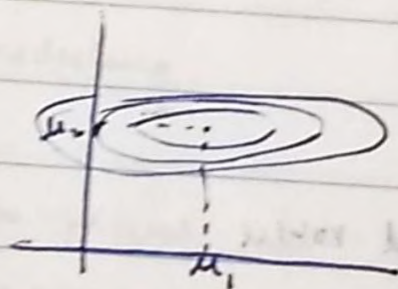
Machine Learning

Independent x_1, x_2

$$x_1 \sim N(\mu_1, \sigma_1)$$

$$x_2 \sim N(\mu_2, \sigma_2)$$

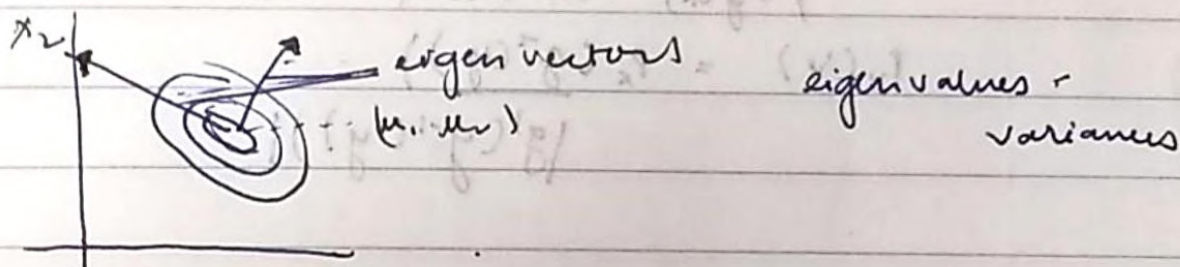
$$p(x_1, x_2) = \frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{1}{2\sigma_1^2}(x_1 - \mu_1)^2} \times \frac{1}{\sqrt{2\pi}\sigma_2} e^{-\frac{1}{2\sigma_2^2}(x_2 - \mu_2)^2}$$



$$p(x_1, x_2) = \left(\frac{1}{\sqrt{2\pi}}\right)^2 \frac{1}{\sigma_1 \sigma_2} e^{-\frac{1}{2} \frac{(x_1 - \mu_1)^2}{\sigma_1^2} - \frac{1}{2} \frac{(x_2 - \mu_2)^2}{\sigma_2^2}}$$
$$= \left(\frac{1}{\sqrt{2\pi}}\right)^2 \frac{1}{\sigma_1 \sigma_2} e^{-\frac{1}{2} \begin{pmatrix} x_1 - \mu_1 & x_2 - \mu_2 \end{pmatrix} \begin{pmatrix} \frac{1}{\sigma_1^2} & 0 \\ 0 & \frac{1}{\sigma_2^2} \end{pmatrix} \begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{pmatrix}}$$

$$P(x) = \left(\frac{1}{2\pi}\right)^{m/2} \frac{1}{\sigma_1 \dots \sigma_m} e^{-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)}$$

Direction and width of ellipse is determined by the eigen decomposition of Σ



$$\Sigma = V \Lambda V^T$$

$$y = V^T (x - \mu)$$

① $X \in \{0, 1\}$ coin

$$Y = 3 + X$$

$$P(X=1) = 0.7$$

$$P(X=0) = 0.3$$

$$Y \in \{3, 4\}$$

$$P(Y=4) = 0.7$$

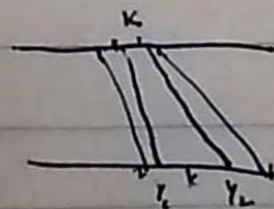
$$P(Y=3) = 0.3$$

② $X \in [1, 2]$ $Y = 2X$

$$P_x(x) = 1$$

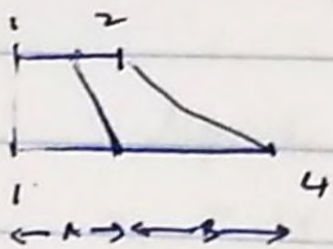
$$\int_0^1 1 dx = x \Big|_0^1 = 1$$

$P_x(x)$



$$y_1 = x_1 \quad y_2 = x_2$$

③ $X \in [1, 2]$ $P_X(x) = 1$ $Y = X^2$



$$\sum_A P_X(x) > \sum_B P_Y(y)$$

$Y = g(X)$ one one & invertible.

$$P_Y(y) = \frac{P_X(g^{-1}(y))}{|g'(g^{-1}(y))|}$$

change of variable for variables in 1D

④ $P_Y(y) = \frac{P_X(x)}{g'(g^{-1}(y))} = \frac{1}{2}$

slope of transformation at its image at y .

⑤ $P_Y(y) = \frac{1}{2\sqrt{y}}$

$$y = g(x) = x^2$$

$$g'(x) = 2x$$

$$g^{-1}(y) = \sqrt{y}$$

$$g'(g^{-1}(y)) = 2\sqrt{y}$$

$$\textcircled{4} \quad x \in [1, 2]$$

$$y = g(x) = x^2$$

$$y \in [1, 4]$$

$$p_y(y) = \text{uniform } [1, 4]$$

$$= \frac{1}{3}$$

$$p_x(x) = ?$$

$$x = f(y) = \sqrt{y}$$

$$p_x(x) = \frac{p_y(f^{-1}(x))}{|f'(f^{-1}(x))|}$$

$$= \frac{p_y(x^2)}{|f'(x^2)|}$$

$$= \frac{p_y(x^2)}{\left| \frac{-1}{2x} \right|}$$

$$= \frac{1/3}{\frac{1}{2|x|}}$$

$$= \frac{1/3}{\frac{1}{2|x|}}$$

$$= \frac{2}{3} |x| = \frac{2x}{3}$$

$$f^{-1}(x) = x^2$$

$$f'(y) = \frac{-1}{2\sqrt{y}}$$

$$f'(f^{-1}(x))$$

$$= \frac{-1}{2\sqrt{x^2}}$$

$$= \frac{-1}{2x}$$

$$P_x(g^{-1}(y)) \cdot |J_{g^{-1}}(y)| = P_x(g^{-1}(y)) \cdot |J_{g^{-1}}(g^{-1}(y))|$$

$$P_x(g^{-1}(y)) \cdot (g^{-1})'(y) = \frac{P_x(g^{-1}(y))}{|g'(g^{-1}(y))|}$$

Multivariate

$$y = g(x) \quad Y = g(x) \quad x \in \mathbb{R}^k \quad y \in \mathbb{R}^k$$

$$P_Y(Y) = \frac{P_x(g^{-1}(Y))}{|J_g(g^{-1}(Y))|}$$

Jacobians

output

$$\begin{pmatrix} y_1 \\ \vdots \\ y_k \end{pmatrix}$$

input $(x_1 \dots x_k)$

$$J = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_k}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \dots & \frac{\partial y_k}{\partial x_2} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_k} & \dots & \frac{\partial y_k}{\partial x_k} \end{pmatrix}$$

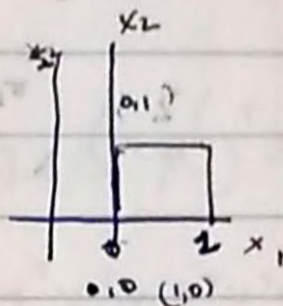
- linear transform.

$$\textcircled{5} Y = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} 2 & 2 \\ 1 & -1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = g(x)$$

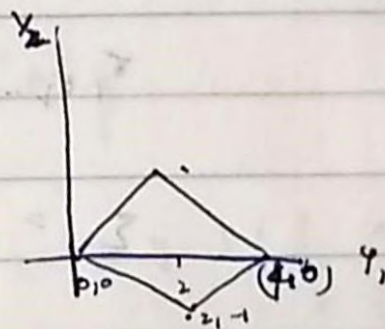
$$X \sim \text{Uniform} [0, 1]$$

$$P_X(X) = 1$$

$$P_Y(Y) = ?$$



$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = g^{-1}(y) = \frac{1}{4} \begin{pmatrix} 1 & 2 \\ 1 & -2 \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}$$



$$J = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} \end{pmatrix} = \begin{pmatrix} 2 & 2 \\ 1 & -1 \end{pmatrix} \quad \text{(Jacobian transform)}$$

$$\text{abs}(|J|) = |-2 - 2| = |-4| = 4$$

$$P_Y(Y) = \frac{P_X(g^{-1}(Y))}{|J(g^{-1}(Y))|} = \frac{1}{4}$$

is diagonal.

$$\Sigma = \begin{pmatrix} \sigma_1^2 & & & \\ & \sigma_2^2 & & \\ & & \dots & \\ & & & \sigma_n^2 \end{pmatrix}$$

$$|\Sigma| = \prod_{i=1}^m \sigma_i$$

$$\sigma_1 \sigma_2 \dots \sigma_m = |\Sigma|^{1/2}$$

$$\Sigma = V \Lambda V^T$$

Capital Lambda.

special linear transform

$$Y = V^T (X - \mu) = g(X)$$

$$X = g^{-1}(Y) = VY + \mu$$

orthonormal.
 $V = (V^T)^{-1}$
 or $V^{-1} = V^T$

$$P_Y(y) = \frac{P_X(g^{-1}(y))}{|J|}$$

$$\Sigma = V \Lambda V^T$$

$$\Sigma^{-1} = V \Lambda^{-1} V^T$$

(see Linear Algebra notes)

$$p_x(x) = (2\pi)^{-m/2} \prod_{i=1}^m |\Sigma|^{-1/2} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

$$p_x(g^{-1}(y)) = (2\pi)^{-m/2} |\Sigma|^{-1/2} e^{-\frac{1}{2}(Vy + \mu - \mu)^T V \Lambda^{-1} V^T (Vy + \mu - \mu)}$$

$$p_x(g^{-1}(y)) = (2\pi)^{-m/2} |\Sigma|^{-1/2} e^{-\frac{1}{2}(Vy) \Lambda^{-1} y^T V^T V \Lambda^{-1} V^T Vy}$$

$$\text{exponent} = -\frac{1}{2} y^T \Lambda^{-1} y$$

var of y is the
eigenvalues of Σ

$$\Rightarrow \frac{1}{2} (y_1, y_2) \begin{pmatrix} \frac{1}{\lambda_1} & 0 \\ 0 & \frac{1}{\lambda_2} \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}$$

$$E[x] = \int p_x(x) x dx$$

$$z = x - \mu = Vy$$

$$dz = dx$$

$$E[x] = \int \left(\frac{1}{2\pi} \right)^{m/2} \frac{1}{|\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)} dx$$

$$= \int e^{-\frac{1}{2} z^T z^{-1} z} (z + \mu) dz$$

symmetric \times asymmetric $\Rightarrow 0$

$$\int \mu dz = \mu \int dz = \mu \cdot \infty$$

$$\int \text{symmetric} \times \text{asymmetric} \Rightarrow 0$$

$$Z = \int \beta(u) g(u) du$$

$$u = -z$$

$$du = -dz$$

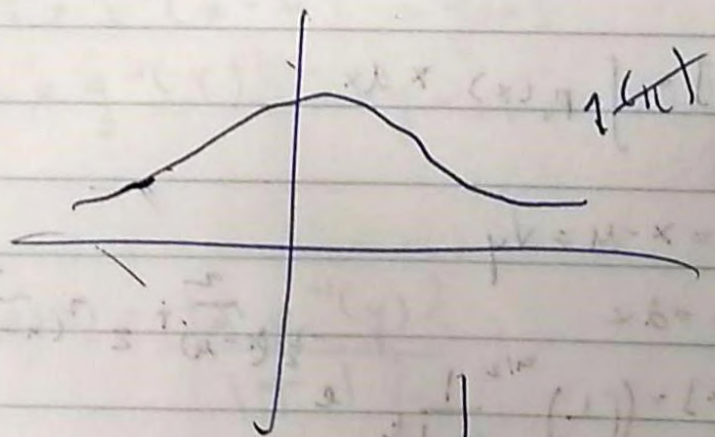
$$I = \int \beta(-z) g(-z) dz$$

$$= \int \beta(z) g(-z) dz$$

$$= \int \beta(z) g(z) dz$$

$$2I = \int \beta(u) g(u) du - \int \beta(u) g(u) du = 0$$

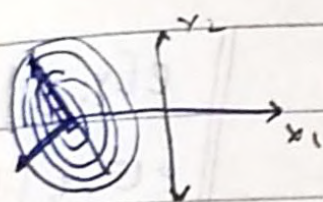
$$g(-u) = -g(u)$$



~~$$g(u) = -g$$~~
~~$$g(-u) = g(u)$$~~

Machine Learning.

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



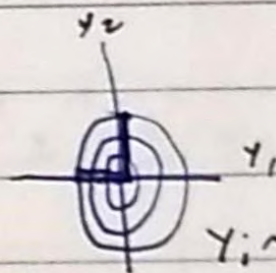
$$Y = V^T (x - \mu)$$

$$X = VY + \mu \Leftrightarrow$$

where

$$\Sigma = V \Lambda V^T$$

$$\Lambda = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_d \end{pmatrix}$$



$$Y_i \sim N(0, \lambda_i)$$

Y_i 's are independent

mean $\mathbb{E}(Y_i) = 0$

Variance λ_i

$$E[x] = \int x \frac{1}{(2\pi)^{d/2}} \frac{1}{|\Sigma|^{1/2}} e^{-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)} dx$$

let $x - \mu = z \Rightarrow x = z + \mu$

$$= \int \frac{1}{(2\pi)^{d/2}} \frac{1}{|\Sigma|^{1/2}} e^{-\frac{1}{2} z^T \Sigma^{-1} z} (z + \mu) dz$$

$$= \int z P(z + \mu) dz + \int \mu P(z + \mu) dz$$

$$\because f(-z) = -f(z)$$

$$1 \times \mu$$

$$I = 0$$

$$\Rightarrow E(x) = \mu$$

$$\text{Cov}(x) = E[(x - \mu)(x - \mu)^T]$$

$$= E[zz^T] = E[VY(VY)^T]$$

$$= E[VY Y^T V^T]$$

$$= V E[YY^T] V^T$$

$$= V \begin{bmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \ddots \\ & & & \lambda_p \end{bmatrix} V^T$$

$$= V \Lambda V^T = \Sigma$$

$$E[y_1 y_3] = ?$$

$$= E[y_1] E[y_3]$$

$$0 \quad 0$$

$$= 0$$

$$E[y_i^2] = E[y_i^2] = ?$$

$$\text{var} = \lambda_i = E[y_i^2] - (E[y_i])^2$$

$$E[y_i] = \lambda_i$$

diagonals

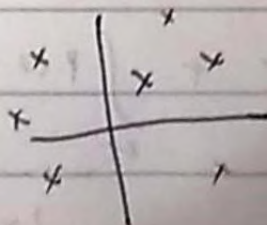
Exercise

independently

Data is generated by a multivariate normal

distribution. Estimate parameters μ, Σ

using Maximum likelihood



data: ~~z~~, $z^{(1)}, z^{(2)}, \dots, z^{(N)}$

Likelihood =

$$\pi(z; \mu, \Sigma) = \prod_{i=1}^N \frac{1}{|\Sigma|^{1/2}} e^{-\frac{1}{2} (z^i - \mu)^T \Sigma^{-1} (z^i - \mu)}$$

$$\log L = \sum_i \left[-\frac{d}{2} \log 2\pi - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (z^i - \mu)^T \Sigma^{-1} (z^i - \mu) \right]$$

$$\frac{\partial a^T B a}{\partial a} = ? \quad \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \quad B \text{ symmetric}$$

$$a^T B a = \sum_i \sum_j a_i a_j B_{ij}$$

$$\frac{\partial (a^T B a)}{\partial a_k} = \sum_j a_j B_{kj} + \sum_i a_i B_{ik}$$

$$= \sum_j a_j B_{jk} + \sum_i a_i B_{ik}$$

$$= \sum_j B_{kj} a_j + \sum_i B_{ki} a_i$$

$$= 2 (k^{\text{th}} \text{ row of } B) a$$

$$\frac{\partial (a^T B a)}{\partial a_k} = 2 b_k^T a$$

$\frac{\partial}{\partial a_k}$

$b_k = k^{\text{th}} \text{ col of } B.$

$$\frac{\partial (a^T B a)}{\partial a} = 2 B a.$$

$\frac{\partial}{\partial a}$

$$\frac{\partial}{\partial \mu} \left(\frac{-1}{2} \sum_i (z^i - \mu)^T \Sigma^{-1} (z^i - \mu) \right) \quad \frac{\partial (a^T b)}{\partial a} = 2b$$

$$\frac{\partial}{\partial \mu} \left[\frac{-1}{2} \left(\sum_i z^i \Sigma^{-1} z^i - 2\mu^T \Sigma^{-1} \sum_i z^i + \mu^T \Sigma^{-1} \mu \right) \right]$$

$$= \begin{pmatrix} 0 \\ -\sum_i z^i \Sigma^{-1} + \Sigma^{-1} \mu \end{pmatrix} = 0$$

$$-\sum_i z^i + \Sigma \mu = 0$$

$$\mu = \frac{\sum_i z^i}{N}$$

$$\frac{\partial \log |A|}{\partial A} = (A^{-1})^T$$

$$\frac{\partial (a^T A^{-1} b)}{\partial A} = -A^{-T} a b^T A^{-1}$$

$$\log L = \sum_i \underbrace{-\frac{1}{2} \log 2\pi}_{\text{constant}} - \frac{1}{2} \sum_i \log |\Sigma| - \frac{1}{2} \sum_i (z^i - \mu)^T \Sigma^{-1} (z^i - \mu)$$

$$\frac{\partial \log L}{\partial \Sigma} = \frac{1}{2} \cdot N (\Sigma^{-1})^T$$

~~constant~~

$$-\frac{1}{2} \sum_i (z^i - \mu) (\Sigma^{-1})^T (z^i - \mu)$$

w

$$\frac{d \log L}{d \Sigma} = \frac{1}{2} (\Sigma^{-1})^T + \frac{1}{2} \Sigma (\Sigma^{-1})^T (Z^i - \mu) (Z^i - \mu)^T (\Sigma^{-1})^T = 0$$

multiply on right by Σ^T

$$\frac{N}{2} \Sigma = \frac{1}{2} \Sigma (\Sigma^{-1})^T \sum_i (Z^i - \mu) (Z^i - \mu)^T$$

$$\Sigma^T = \frac{1}{N} \sum_i \underbrace{(Z^i - \mu) (Z^i - \mu)^T}_{\text{Symmetric}}$$

$$\Rightarrow \Sigma^T = \Sigma = \frac{1}{N} \sum_i (Z^i - \mu) (Z^i - \mu)^T$$

Develop derivative identities

ex. 1

$\frac{d}{dx} (AB)_{ij}$
 ← matrices
 ← scalar

$$C = AB$$

$$C_{ij} = \sum_k A_{ik} B_{kj}$$

$$\left(\begin{array}{c} i \\ \hline \end{array} \right) \left(\begin{array}{c} | \\ j \end{array} \right)$$

$$\frac{d}{dx} C_{ij} = \sum_k \frac{\partial A_{ik}}{\partial x} B_{kj} + \frac{\partial B_{kj}}{\partial x} A_{ik}$$

$$\left(\frac{\partial C}{\partial x} \right)_{ij} = \left(\frac{\partial A}{\partial x} B \right)_{ij} + \left(A \frac{\partial B}{\partial x} \right)_{ij}$$

$$\frac{\partial C}{\partial x} = \frac{\partial A}{\partial x} B + A \frac{\partial B}{\partial x} \quad (\text{matrix product rule})$$

$$0 = \frac{\partial}{\partial x} (A^{-1}) (A^{-1})^T (I) = \frac{\partial}{\partial x} (A^{-1} A^{-1}^T) (I)$$

$$\frac{\partial}{\partial x} (A^{-1} A) = \frac{\partial}{\partial x} (I)$$

$$= \frac{\partial A^{-1}}{\partial x} A + A^{-1} \frac{\partial A}{\partial x} = 0$$

$$\frac{\partial A^{-1}}{\partial x} = -A^{-1} \frac{\partial A}{\partial x} A^{-1}$$

ex. 3

$$\frac{\partial A^{-1}}{\partial x_{kl}} = -A^{-1} \frac{\partial A}{\partial x_{kl}} A^{-1}$$

$$= A^{-1} \begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix} A^{-1}$$

$$\frac{\partial}{\partial d} \begin{pmatrix} a & b & c \\ d & e & f \\ g & h & i \end{pmatrix} = \begin{pmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}$$

ex. 4

$$\frac{\partial a^T A^{-1} b}{\partial x_{kl}} \quad \left. \begin{array}{l} \text{scalar} \\ \text{scalar} \end{array} \right\}$$

$$= a^T \frac{\partial A^{-1}}{\partial x_{kl}} b$$

$$= -a^T A^{-1} \begin{pmatrix} \dots & 0 \\ \dots & 1 & \dots \\ \dots & \dots & 0 \end{pmatrix} A^{-1} b$$

$$A^T \begin{pmatrix} \text{all zeros} \\ \text{except } k\text{-th} \\ \text{which is } 1 \end{pmatrix} = \begin{pmatrix} \text{zeros except} \\ \text{the } k\text{-th column} \end{pmatrix}$$

The diagram shows a column vector with a '1' in the \$k\$-th position and zeros elsewhere. To its right is a matrix with a single column highlighted, representing the \$k\$-th column of the matrix \$A\$.

k^{th} column of A^{-1} moves to k^{th} column of result

$$\begin{pmatrix} 1 \\ \vdots \\ 0 \end{pmatrix} \begin{pmatrix} A^{-1} \end{pmatrix} = \begin{pmatrix} i \\ \vdots \\ 1 \end{pmatrix} A^{-1}_{ik} A^{-1}_{kj}$$

\hat{A} = Matrix whose ij entry is $A^{-1}_{ik} A^{-1}_{kj}$

$$\frac{\partial a^T A^{-1} b}{\partial A} = -a^T \hat{A} b$$

$$\frac{\partial A}{\partial A}$$

$$= -\sum_i \sum_j a_i b_j \hat{A}_{ij}$$

$$= -\sum_i \sum_j a_i b_j A^{-1}_{ik} A^{-1}_{kj}$$

$$\Rightarrow \frac{\partial a^T A^{-1} b}{\partial A} = \mathbf{f}_A$$

where entries in \hat{A} are

$$\frac{\partial (a^T A^{-1} b)}{\partial A} = -A^{-1T} a b^T A^{-1}$$

Linear Regression

$$L = p(t|w)$$

$$t|w \sim N(\phi w, \frac{1}{\beta} \Sigma)$$

$$t_i|w \sim N(w^T \phi(x_i), \frac{1}{\beta}) \quad t_i\text{'s are independent}$$

Normal prior for w :

$$w \sim N(m_0, S_0)$$

m_0 is the mean of the prior

= mean after seeing 0 observations

S_0 is the covariance matrix of prior

cov. after seeing 0 observations.

$$\rightarrow w|t \propto p(t|w) p(w)$$

$$\frac{1}{e} (t - \phi w)^T \left(\frac{1}{\beta} \Sigma\right)^{-1} (t - \phi w) \quad \frac{1}{e} (w - m_0)^T S_0^{-1} (w - m_0)$$

multiply & see what this looks like in terms of w

Machine Learning

Linear Regression

$$t_i \sim N(w^T \phi(x_i), \frac{1}{\beta})$$

t_i independent

$$t|w \sim N(\Phi w, \frac{1}{\beta} I)$$

$$\left\| \begin{bmatrix} \frac{1}{\beta} \\ \frac{1}{\beta} \\ \dots \end{bmatrix} \right\|$$

$$w \sim N(m_0, S_0) \quad (\text{gaussian prior})$$

$$x \sim N(\mu_x, S_x)$$

$$y|x \sim N(Ax + b, S_y)$$

observe y & want $P(w | y = y)$

$$\begin{array}{l} t \equiv y \\ w \equiv x \end{array}$$

Generic MVN, $z \sim N(\mu, \Sigma)$

$$\propto \left(\frac{1}{2\pi}\right)^{d/2} \frac{1}{|\Sigma|^{1/2}} e^{-\frac{1}{2}(z-\mu)^T \Sigma^{-1} (z-\mu)}$$

$$e^{-\frac{1}{2} z^T \Sigma^{-1} z} \underbrace{e^{-\frac{1}{2} \mu^T \Sigma^{-1} \mu}}_{\text{does not depend on } z} e^{z^T \Sigma^{-1} \mu}$$

does not depend on z
 \rightarrow ignore

$$P(z) \propto e^{-\frac{1}{2} z^T \Sigma^{-1} z} e^{z^T \Sigma^{-1} \mu}$$

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

$$= \left(e^{-\frac{1}{2} x^T S_x x} \quad e^{-x^T S_x^{-1} m_x} \right) \left(e^{-\frac{1}{2} (y - (Ax+b))^T S_y^{-1} (y - (Ax+b))} \right)$$

consider exponents of second term

$$-\frac{1}{2} y^T S_y^{-1} y - \frac{1}{2} (Ax+b)^T S_y^{-1} (Ax+b) + (Ax+b)^T S_y^{-1} y$$

$$\textcircled{3} -\frac{1}{2} x^T A^T S_y^{-1} A x$$

$$-x^T A^T S_y^{-1} b$$

$$\textcircled{4} -\frac{1}{2} b^T S_y^{-1} b$$

$$\textcircled{5} x^T A^T S_y^{-1} y + b^T S_y^{-1} y$$

① & ③ quadratic terms

$$-\frac{1}{2} x^T (S_x^{-1} + A^T S_y^{-1} A) x \Rightarrow \Sigma^{-1} = S_x^{-1} + A^T S_y^{-1} A$$

on comparison

②, ④ & ⑤

$$x^T (S_x^{-1} m_x - A^T S_y^{-1} b + A^T S_y^{-1} y)$$

$$= x^T (S_x^{-1} m_x + A^T S_y^{-1} (y - b))$$

$$\Rightarrow \Sigma^{-1} \mu = S_x^{-1} m_x + A^T S_y^{-1} (y - b)$$

on comparison

$$\Rightarrow \mu = (S_x^{-1} + A^T S_y^{-1} A)^{-1} (S_x^{-1} m_x + A^T S_y^{-1} (y - b))$$

$$\left(\frac{1}{\beta} I\right)^T = \beta I$$

Posterior

$$w|t \sim N(\mu^*, \Sigma^*)$$

$$m_y = m_0, S_y = S_0$$

$$K = \Phi, b = 0$$

$$S_y = \frac{1}{\beta} I$$

$$\Sigma^* = \left(S_0^{-1} + \Phi^T \left(\frac{1}{\beta} I\right) \Phi\right)^{-1}$$

$$\begin{aligned} \mu^* &= \Sigma^* \left(S_0^{-1} m_0 + \Phi^T \beta I \bar{t}\right) \\ &= \Sigma^* \left(S_0^{-1} m_0 + \beta \Phi^T \bar{t}\right) \end{aligned}$$

MAP = ~~maximum~~ input that maximizes posterior
(in normal dist., it is same as mean)
= μ^* .

$$w|t \sim N(\mu_N, S_N)$$

$$S_N = \left(S^{-1} + \beta \Phi^T \Phi\right)^{-1}$$

$$\mu_N = S_N \left(\beta \Phi^T t + S_0^{-1} m_0\right)$$

'no information prior' (fall back to MLE)

$$\text{variance} \rightarrow \infty \quad S_0 = \begin{pmatrix} \infty & & \\ & \infty & \\ & & \ddots \end{pmatrix} \quad S_0^{-1} = \begin{pmatrix} 0 & & \\ & 0 & \\ & & \ddots \end{pmatrix}$$

$$S_N = \left(\beta \Phi^T \Phi\right)^{-1}$$

$$\mu_N = \left(\beta \Phi^T \Phi\right)^{-1} \left(\beta \Phi^T t + 0\right)$$

$$= \left(\beta \Phi^T \Phi\right)^{-1} \left(\beta \Phi^T t\right)$$

$$= \frac{1}{\beta} \left(\Phi^T \Phi\right)^{-1} \beta \left(\Phi^T t\right) = \left(\Phi^T \Phi\right)^{-1} \left(\Phi^T t\right)$$

$$w \sim N(0, \frac{1}{\alpha} I)$$

α
alpha

$$S_N = \left(\frac{1}{\alpha} I + \beta \Phi^T \Phi \right)^{-1}$$

$$w_N = \left(\frac{1}{\alpha} I + \beta \Phi^T \Phi \right)^{-1} (\beta \Phi^T t + 0)$$

$$= \left(\frac{1}{\alpha} I + \Phi^T \Phi \right)^{-1} (\Phi^T t)$$

$$= \left(\lambda I + \Phi^T \Phi \right)^{-1} (\Phi^T t)$$

$$\text{let } \lambda = \frac{\alpha}{\beta}$$

$$w_N = (\lambda I + \Phi^T \Phi)^{-1} (\Phi^T t)$$

regularized linear regression

Predictive distribution

$$P(t_{N+1} | w) \sim N(\phi(x_{N+1})^T w, \frac{1}{\beta})$$

$$P(t_{N+1} | t) \sim ?$$

$$w \sim N(\mu_N, \Sigma_N)$$

$$t_{n+1} | w \sim N(\phi(x_{n+1})^T w, \frac{1}{\beta})$$

$$a = \phi^T$$

$$b = 0$$

$$t_{n+1} \sim N(\phi(x_{n+1})^T \mu_N,$$

$$\frac{1}{\beta} + \phi(x_{n+1})^T \Sigma_N \phi(x_{n+1}))$$

What if β is not known?

conjugate prior/posterior for (w, β)

Likelihood

$$\prod_{n=1}^N \frac{1}{\sqrt{2\pi}} \frac{1}{\beta} e^{-\frac{1}{2\beta} (\phi w - t)^T (\phi w - t)}$$

$$\left(\frac{1}{2\pi} \right)^{N/2}$$

$$\beta^{-N} e^{-\frac{1}{2\beta} (\phi w - t)^T (\phi w - t)}$$

$$\left\| \begin{array}{l} \Sigma = \frac{1}{\beta} \Sigma \\ |Z| \propto \left(\frac{1}{\beta}\right)^N \\ \frac{1}{|\Sigma|} = \beta^{N/2} \end{array} \right.$$

What would be a conjugate prior for β ?

$$\beta \propto \beta^{-p} e^{-\beta}$$

$$\left\| \begin{array}{l} \text{Gamma dist.} \\ \text{Gamma}(x | a, b) \\ = \frac{1}{\Gamma(a)} b^a x^{a-1} e^{-bx} \end{array} \right.$$

(ay)

$$X \begin{cases} w \sim N(\mu_0, \Sigma_0) \\ \beta \sim \text{Gamma}(\beta, a, b) \end{cases} \quad \left. \begin{array}{l} \text{independent} \\ \text{independent} \end{array} \right\}$$

$$p(\beta, w | G)$$

$$S_w = (\Sigma_0^{-1} + \beta \phi^T \phi)^{-1} \quad \text{dependent}$$

$$(1000) \quad \text{unstable} \quad S_w = \frac{1}{\beta} \Sigma_0$$

Machine Learning

Goal: use bayesian solution for both w, β .

Problem: $\left\{ \begin{array}{l} \text{even if we start with independent } P(w) \& P(\beta) \\ \text{in prior,} \\ \text{posterior will have } P(w, \beta) \text{ with } w, \beta \text{ correlated} \end{array} \right.$

$$t_i \sim N(w^T \phi(x_i), \frac{1}{\beta})$$

$$t \sim N(\Phi w, \frac{1}{\beta} I)$$

Gamma Distribution

$$P(x|a, b) = \frac{1}{\Gamma(a)} b^a x^{a-1} e^{-bx} \quad a, b, x > 0$$

Distribution over positive random variable

Prior

$$P(\beta | a_0, b_0) = \text{Gamma}(\beta, a_0, b_0) \quad \dots \text{marginal}$$

$$P(w | m_0, s_0, \beta) = N(m_0, \frac{1}{\beta} s_0) \quad \dots \text{conditional}$$

Hope for posterior

$$P(\beta | a_N, b_N) = \text{Gamma}(\beta | a_N, b_N)$$

$$P(w | m_N, s_N, \beta) = N(m_N, \frac{1}{\beta} s_N)$$

Posterior \propto Prior \times Likelihood

$$\beta \text{ prior} \dots \frac{1}{\Gamma(a_0)} b_0^{a_0} \beta^{a_0-1} e^{-b_0 \beta}$$

$$w \text{ prior} \dots \left(\frac{1}{2\pi}\right)^{d/2} \frac{1}{|\frac{1}{\beta} S_0|^{1/2}} e^{-\frac{1}{2} (w-m_0)^T (\frac{1}{\beta} S_0)^{-1} (w-m_0)}$$

$$\text{Likelihood} \dots \left(\frac{1}{2\pi}\right)^{N/2} \frac{1}{\sqrt{|\frac{1}{\beta} I|}} e^{-\frac{1}{2} (t-\Phi w)^T \beta I (t-\Phi w)}$$

$$\sqrt{\left(\frac{1}{\beta}\right)^N} = \left(\frac{1}{\beta}\right)^{N/2}$$

Param	(w, β)	prior (w, β)
data	t	
posterior	$p(w t)$	

$$\left(\frac{1}{\beta} S_0\right)^{1/2} = \left(\frac{1}{\beta}\right)^{d/2} |S_0|^{1/2} \Rightarrow \frac{1}{\sqrt{|\frac{1}{\beta} S_0|}} = \beta^{d/2} \frac{1}{\sqrt{|S_0|}}$$

ignoring constant terms

$$\beta^{a_0-1} e^{-b_0 \beta} e^{-\frac{\beta}{2} w^T S_0^{-1} w} e^{-\frac{\beta}{2} m_0^T S_0^{-1} m_0} \beta w^T S_0^{-1} m_0$$

$$\beta^{N/2} e^{-\frac{\beta}{2} t^T t} e^{-\frac{\beta}{2} w^T \Phi^T \Phi w} \beta w^T \Phi^T t$$

Quadratic term

$$-\frac{\beta}{2} w^T (\beta_0^{-1} + \Phi^T \Phi) w \quad \text{vls} \quad -\frac{1}{2} z^T \Sigma^{-1} z$$

$$\Rightarrow \Sigma^{-1} = \beta_0^{-1} + \Phi^T \Phi$$

$$\Sigma = \frac{1}{\beta} \underbrace{(\beta_0^{-1} + \Phi^T \Phi)^{-1}}_{S_N} \quad S_N = (\beta_0^{-1} + \Phi^T \Phi)^{-1}$$

Linear term

$$w^T (\beta_0^{-1} m_0 + \beta \Phi^T t) \quad \text{vls} \quad z^T \Sigma^{-1} \mu$$

$$\boxed{w} = \Sigma^{-1} \mu$$

$$\Sigma^{-1} \mu = \boxed{}$$

$$\mu = \Sigma \boxed{}$$

$$\mu = \Sigma(\beta) (\beta_0^{-1} m_0 + \Phi^T t)$$

$$\mu = \frac{1}{\beta} (\beta_0^{-1} + \Phi^T \Phi)^{-1} \beta (\beta_0^{-1} m_0 + \Phi^T t)$$

$$\boxed{\mu = (\beta_0^{-1} + \Phi^T \Phi)^{-1} (\beta_0^{-1} m_0 + \Phi^T t)}$$

next work on β

but

look at posterior on w .

$$\frac{1}{(2\pi)^{d/2}} \frac{1}{\sqrt{\frac{1}{\beta} S_N}} e^{-\frac{1}{2} (w^T (\frac{1}{\beta} S_N)^{-1} w - \frac{1}{2} m_N^T (\frac{1}{\beta} S_N)^{-1} m_N)} \frac{1}{\beta} \frac{1}{\sqrt{|S_0|}} e^{-\frac{1}{2} m_N^T (\frac{1}{\beta} S_N)^{-1} m_N}$$

not accounted for

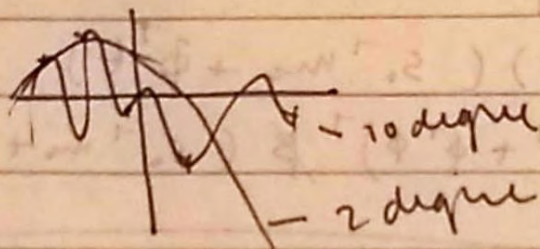
$$\text{Gamma}(x | a, b) \propto x^{a-1} e^{-bx}$$

$$\beta^{a_0 + \frac{N}{2} - 1} \Rightarrow a_N = a_0 + \frac{N}{2}$$

$$e^{-\frac{1}{2} \left(b_0 + \frac{1}{2} m_0^T S_0^{-1} m_0 + \frac{1}{2} t^T t - \frac{1}{2} m_N^T S_0^{-1} m_N \right)}$$

β_N

\therefore no term remains, by completing the squares the form of the posterior is Normal \times Gamma.



$$\Phi = \begin{pmatrix} 1 & x & x^2 & \dots & x^{10} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \vdots & \ddots & \vdots \end{pmatrix}$$

$$\Phi = \begin{pmatrix} 1 & x & x^2 \\ \vdots & \vdots & \vdots \end{pmatrix}$$

$$\begin{pmatrix} 1 \\ x \\ x^2 \\ \vdots \end{pmatrix}$$

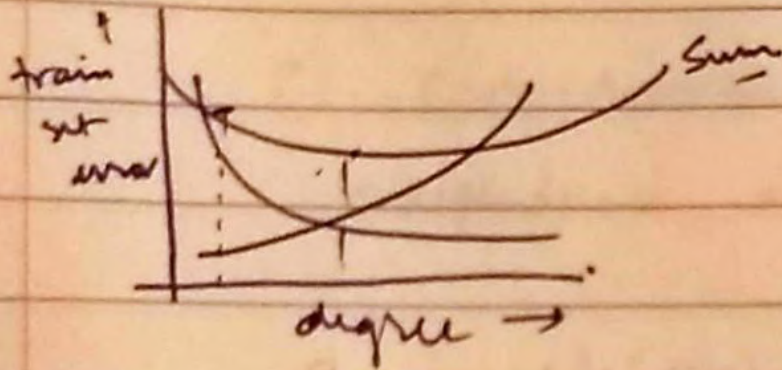
Model Selection

How to pick extra parameters

$\hat{=}$ param of prior

Simple intuition

Bayesian information criteria



Size of confidence interval grows with complexity of model class

Evidence maximization

Machine Learning

Model Selection

$$P(w) = N(0, \frac{1}{\alpha} \Sigma)$$

$$P(t|w) = N(\Phi w, \frac{1}{\beta} \Sigma)$$

Posterior $P(w|t) \sim N(m_N, S_N)$

$$m_N = \beta S_N \Sigma^{-1} t$$

$$S_N = (\alpha \Sigma + \beta \Phi^T \Phi)^{-1}$$

what values of α, β ?

hyper-parameters

polynomial regression

$$\phi(x) = \begin{pmatrix} 1 \\ x \\ \dots \\ x^d \end{pmatrix}$$

what value of d ?

Problem 1

fixed d

pick α, β

So as to get "best" posterior

Problem 2

pick α, β, d

!

Problem

θ parameters

γ hyperparameters

prior $P(\theta|\gamma)$

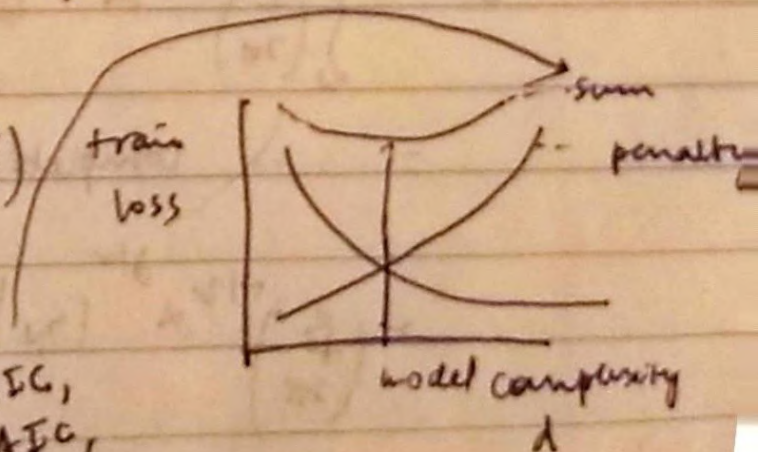
Likelihood $P(t|\theta)$

find $P(\theta|t, \gamma)$

what value for γ

Standard trick

min $Loss(\gamma) + Penalty(\gamma)$
to pick γ .



BIC, AIC, MDL, SRM

Pick γ using evidence function
 Evidence $P(t | Y)$
 \equiv Max. likelihood on hyper parameters
 \equiv 2nd level max. L
 \equiv Empirical Bayes
 \equiv Evidence approximation.

Bishop: ① intuitive arguments that this works
 ② rough calculation shows that we get BIC as a special case.

$$\text{Evidence}(Y) = \int_{\theta} \underbrace{P(\theta | Y)}_{\text{prior}} \underbrace{P(t | \theta)}_L d\theta$$

pick hyper param to max

Model

$$\hat{E}_V \approx \int_{\omega} \left(\frac{L}{2\pi} \right)^{d/2} \alpha^{d/2} e^{-\frac{\alpha}{2} \omega^T \omega} \left(\frac{L}{2\pi} \right)^{N/2} \beta^{N/2} e^{-\frac{\beta}{2} (\phi \omega - t)^T (\phi \omega - t)}$$

\approx .. completing the square to normal

$$\approx \left(\frac{\beta}{2\pi} \right)^{N/2} \alpha^{N/2} |\Sigma| e^{-\frac{\beta}{2} u^T \Sigma^{-1} u - t^T u} e^{-\frac{\alpha}{2} u^T \Sigma^{-1} u}$$

1. Calculate Log Likelihood
2. Take derivative
3. Solve for λ, β .

Solution gives an iterative algorithm:

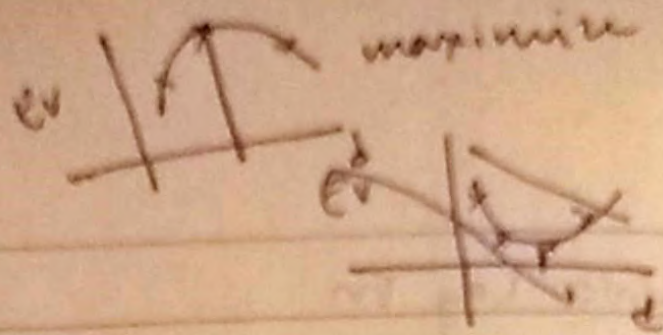
1. Calc. λ_i - eigenvalues of $\beta \phi^T \phi$
2. Calc. $\gamma = \sum_{i=1}^d \frac{\lambda_i}{\lambda_i + \beta}$

3. Update $\alpha = \frac{\gamma}{\|\mathbf{m}_N\|^2}$

$$\frac{1}{\beta} = \frac{1}{N-1} \|\phi \mathbf{m}_N - \mathbf{t}\|^2$$

Model Selection Algorithm

1. Initialize λ, β
2. Repeat until convergence
 - Calculate $\mathbf{m}_{N+1}, \mathbf{s}_N$
 - Update λ, β



Model 2,

alg. 10 Select α, β, d
for $d = 1, \dots, D$

run alg. for Model Selection #1 ^{to pick α, β}

Calculate evidence using d, α, β

Pick d, α_d, β_d which ^{max} Evidence.

eigenvalues

$$\Phi^T \Phi \quad \hat{\lambda}_i$$

$$\beta \Phi^T \bar{q} \quad \lambda_i = \beta \hat{\lambda}_i$$

$$\alpha \Gamma + \beta \bar{q}^T \Phi \quad \alpha + \beta \hat{\lambda}_i$$

$$S_N \quad \frac{1}{\lambda_i + \alpha} = \frac{1}{\beta \hat{\lambda}_i + \alpha}$$

$$\lambda_i \rightarrow \infty \Rightarrow \text{var in direction} \approx 0$$

$$\lambda_i \rightarrow 0 \Rightarrow \text{var} \approx \infty$$

$$\frac{\lambda_i}{\lambda_i + \alpha} \begin{matrix} (\lambda_i \rightarrow \infty) \rightarrow 1 \\ (\lambda_i \rightarrow 0) \rightarrow 0 \end{matrix}$$

$$r = \sum \frac{\lambda_i}{\lambda_i + \alpha} \approx \text{no. of determined dimensions}$$

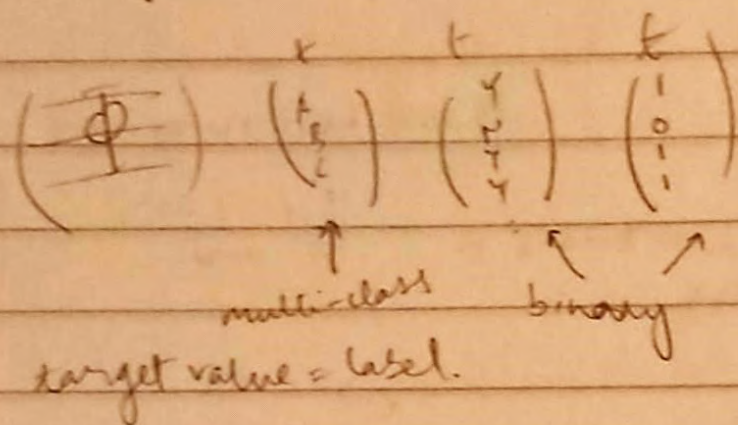
Machine Learning

Linear Regression

limitations - model is linear

predicting $t \in \mathbb{R}$

Classification: when t is discrete



Simplest Soln for 2 class

Use linear regression

$$Y \leftrightarrow 1$$

$$N \leftrightarrow -1$$

use $\max L \rightarrow W$

for new example x ,

compute $a = w^T x$

if $a \geq 0 \Rightarrow$ say yes
otherwise \Rightarrow no

Methodology

Specify a model that explains how data is generated

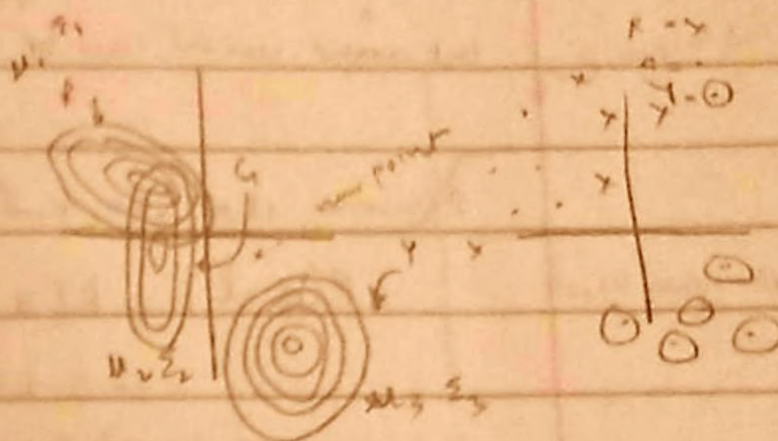
Then given data learn parameters of that model. (or a posterior over parameters)

Example

$$P(c=1) \quad P(P) = 0.5$$

$$P(c=2) \quad P(G) = 0.4$$

$$P(c=3) \quad P(Y) = 0.1$$



$$P(x|P)$$

$$P(x|G)$$

$$P(x|Y)$$

To generate data

for each i

pick class $c_i \in \{1, 2, 3\}$

pick x from $P(x|c = c_i)$

We also want to include cases where features are discrete

$$x = \begin{pmatrix} 0 & 1 & 1 \\ 0 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 1 \end{pmatrix}$$

We will need a different $P(x|C=l)$
but some model can be used

Assume that we know

$$P(C=l) \quad P(x|C=l)$$

How is a new example classified?

compute $P(C=j|x)$

$$P(C=j|x) = \frac{P(x|C=j) P(C=j)}{P(x)}$$

$$= \frac{P(C=j) P(x|C=j)}{\sum_i P(x|C=i) P(C=i)}$$

$$= \frac{P(C=j) P(x|C=j)}{\sum_i P(x|C=i) P(C=i)}$$

$$\text{Define } a_j = \log [P(C=j) P(x|C=j)]$$

$$e^{a_j}$$

$$\sum_i e^{a_i}$$

← Softmax

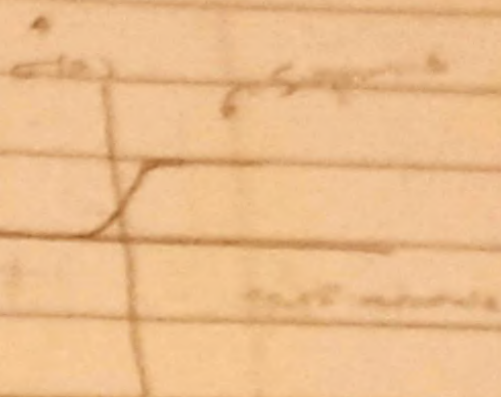
In 2 class case

$$P(c=1) = \frac{e^{a_1}}{e^{a_1} + e^{a_2}}$$

$$= \frac{1}{1 + e^{a_2 - a_1}}$$

Define a to be $a_1 - a_2$

$$\Rightarrow \frac{1}{1 + e^{-a}}$$



with equal costs, predict $c=1$

$$\Leftrightarrow P(c=1) \geq 1/2$$

$$\Leftrightarrow a \geq 0$$

What does prediction look like in geometric space
when $P(y|c_j) = N(\mu_j, \Sigma_j)$

Consider cost with $\mu=2$, $\Sigma_1 = \Sigma_2 = \Sigma$

When do we predict class = 1

$$a = a_1 - a_2$$

$$= \log \frac{P(c=1) P(x|c=1)}{P(c=2) P(x|c=2)}$$



$$\begin{aligned}
 a &= \log P(c=1) - \log P(c=2) \\
 &= \frac{d}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x - \mu_1)^T \Sigma^{-1} (x - \mu_1) \\
 &\quad + \frac{d}{2} \log(2\pi) + \frac{1}{2} \log |\Sigma| + \frac{1}{2} (x - \mu_2)^T \Sigma^{-1} (x - \mu_2) \\
 &\geq \frac{1}{2}
 \end{aligned}$$

$$\frac{1}{2} x^T \Sigma^{-1} x - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 \neq x^T \Sigma^{-1} \mu_1 \quad \text{linear in } x$$

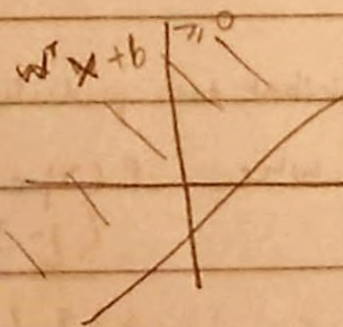
Cancel out

$$\frac{1}{2} x^T \Sigma^{-1} x + \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2 + x^T \Sigma^{-1} \mu_2 \geq \frac{1}{2}$$

$$w = \Sigma^{-1} (\mu_2 - \mu_1)$$

$$b = \frac{1}{2} + \log \left(\frac{P(c=1)}{P(c=2)} \right) - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 - \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2$$

$$x^T w + b \geq 0$$

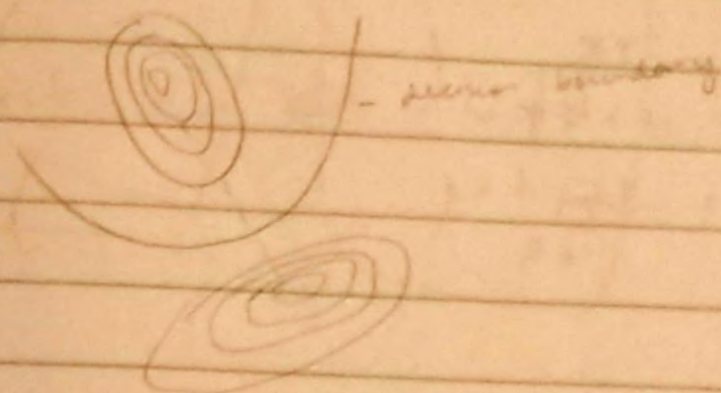


When $\Sigma_1 \neq \Sigma_2$

only one step in the derivation changes

$$\frac{1}{2} x^T \Sigma_2^{-1} x - \frac{1}{2} x^T \Sigma_1^{-1} x = x^T \left(\frac{1}{2} (\Sigma_2^{-1} - \Sigma_1^{-1}) \right) x$$

Quadratic



or

4 classes

Max likelihood estimate

$$\{ \mu_i, \Sigma_i \}$$

$$P_i = \pi(\omega_i)$$

$$\text{Likelihood} = \left[\prod_{i \in A} P_A N(x_i | \mu_A, \Sigma_A) \right]$$

$$\left[\prod_{i \in B} P_B N(x_i | \mu_B, \Sigma_B) \right]$$

$$\left[\prod_{i \in Y} P_Y N(x_i | \mu_Y, \Sigma_Y) \right]$$

$$\left[\prod_{i \in B} P_B N(x_i | \mu_B, \Sigma_B) \right]$$

each μ_j is independent of μ_i .

$$\text{Log } L = \left(\sum_{i \in A} \log P_A + \sum_{i \in G} \log N(x_i | \mu_G, \Sigma_G) \right)$$

$$+ \sum_{i \in B} \left(\dots \right)$$

$$+ \sum_{i \in Y} \left(\dots \right)$$

$$+ \sum_{i \in B} \left(\dots \right)$$

For μ_i, Σ_i same as above L for M VN params

$$\text{For } P_A + P_B + P_Y = 1$$

$$P_A = \frac{\# \text{ points } G}{\text{total \# of points}}$$

$$\bar{\Sigma} = \frac{1}{N} \sum (x_i - \mu) (x_i - \mu)^T$$

Grand covariance

$$\bar{\Sigma} = \frac{1}{N} \left(\sum_G (x_i - \mu_G) (x_i - \mu_G)^T \right.$$

$$+ \sum_B (x_i - \mu_B) (x_i - \mu_B)^T$$

$$\left. + \dots \right)$$

data \approx lot of params

$\{ \mu_i, \sigma_i, \gamma_i \}$

1 + 1 + 1 + d^2 params

param \approx w, b

Predict (a) \odot w^{-1} \approx 15 \approx 0

\approx 3 params

Machine Learning

$$P(c=c)$$
$$P(x|c=c)$$

observed x_1, \dots, x_n

\Rightarrow Max L for param of $P(c=c) P(x|c=c)$

$$\{\mu_i, \Sigma_i\}_{i=1, \dots, K}$$

To predict on new example

$$w = f(\{\mu_i, \Sigma_i\})$$
$$b = g(\dots)$$

$$Z \text{ is yes} \Leftrightarrow w^T \phi(x) + b \geq 0$$

instead of writing $w^T \phi(x) + b$

$$w = (w_1, \dots, w_d)$$
$$v = (\underbrace{w_0}_{\text{representable}}, w_1, \dots, w_d)$$

$$\hat{\phi}(x) = (1, \phi_1(x), \dots, \phi_d(x))$$

$$v^T \hat{\phi}(x) = w_0 + w^T \phi(x)$$

$$\text{set } w_0 = b$$

linear regression

$$a_i = w^T \phi(x_i)$$

$$y_i = a_i$$

$$t_i \sim N(y_i, \frac{1}{\beta})$$

logistic regression

$$a_i = w^T \phi(x_i)$$

$$y_i = \sigma(a_i)$$

Sigmoid

$$t_i \sim \text{Bernoulli}(y_i)$$

drawn independently

random variable

$$\Phi = \begin{pmatrix} \text{---} \\ \text{---} \\ \text{---} \end{pmatrix}$$

$$t = \begin{pmatrix} 0 \\ 1 \\ 0 \\ 1 \end{pmatrix}$$

Prob that
label of ith
example is 1
 $= y_i$

$$\Rightarrow w_{max} L$$

$$L = \prod_i y_i^{t_i} (1-y_i)^{1-t_i}$$

$$\log L = \sum_i t_i \log y_i + (1-t_i) \log(1-y_i)$$

diff wrt w ,

$$\frac{d \log L}{d w} = \sum_i t_i \frac{1}{y_i} \frac{\partial y_i}{\partial w} + (1-t_i) \left(\frac{-1}{1-y_i} \right) \frac{\partial y_i}{\partial w}$$

$$\frac{\partial y_i}{\partial w} = \frac{\partial y_i}{\partial a_i} \frac{\partial a_i}{\partial w}$$

$$y = \sigma(a)$$

$$\frac{dy}{da} = \frac{d}{da} \left(\frac{1}{1+e^{-a}} \right)$$

$$= (-1) \frac{1}{(1+e^{-a})^2} (e^{-a})(-1)$$

$$= \frac{1}{1+e^{-a}} \cdot \frac{e^{-a}}{1+e^{-a}}$$

$$= \sigma(a) (1 - \sigma(a))$$

$$\frac{\partial a_i}{\partial w} = \frac{d}{dw} w^T \phi(x_i) = \phi(x_i)$$

$$\frac{d \log L}{dw} = \sum_i t_i \frac{1}{y_i} \frac{dy_i}{da_i} \frac{\partial a_i}{\partial w} + \sum_i (1-t_i) \frac{1}{1-y_i} \frac{dy_i}{da_i} \frac{\partial a_i}{\partial w} (-1)$$

$$= \sum_i \frac{t_i}{y_i} y_i (1-y_i) \phi(x_i) - \sum_i \left(\frac{1-t_i}{1-y_i} \right) y_i (1-y_i) \phi(x_i)$$

$$= \sum_i \left[t_i \phi(x_i) - t_i y_i \phi(x_i) - y_i \phi(x_i) + t_i y_i \phi(x_i) \right]$$

$$= \sum_i \left[(t_i - y_i) \phi(x_i) \right] = \phi^T (t - y)$$

$$\frac{d \log L}{dw} = 0$$

vector elements matrix columns

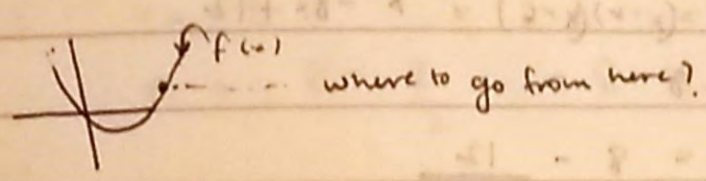
$$\sum_i t_i \phi(x_i) - \sum_i y_i \phi(x_i) = 0$$

log L is concave and has only one maxima. So, GD will converge.

$$\sum_i t_i \Phi(x_i) = \sum_i y_i \Phi(x_i) \\ = \sum \sigma(w^T \Phi(x_i)) \Phi(x_i)$$

No simple "closed form solutions" for w .

optimize directly.



- if f is increasing go left
- if f is decreasing go right

To maximize, go with gradient
To minimize, go against gradient

Gradient Descent:

initialize x

Repeat

$$x \leftarrow x - \eta \cdot f'(x)$$

Gradient Descent for Logistic Regression

$$w \leftarrow w + \eta \underbrace{\Phi^T(Lt - y)}_{\frac{d}{dw} \log L}$$

$$\equiv w \leftarrow w + \eta \frac{d}{dw} \log L$$

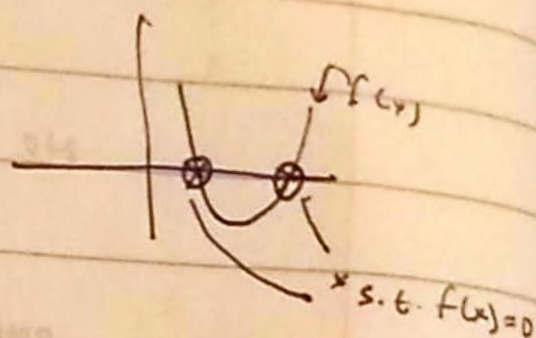
← Maximize log likelihood.

if fn. is linear
this method terminates in
1 step and is exact

Newton's method.

To find a zero of $f(x)$

$$x_{n+1} \leftarrow x_n - \frac{f(x_n)}{f'(x_n)}$$



ex. $f(x) = (x-2)(x-6) = x^2 - 8x + 12$

$$x_0 = 8$$

$$x_1 = 8 - \frac{12}{8}$$

why?

$$f(x_0+h) \approx f(x_0) + f'(x_0)h + \left\{ \frac{1}{2} f''(x_0)h^2 + \dots \right.$$

Taylor
expansion

$$\approx f(x_0) + f'(x_0)h$$

higher order terms
are small; ignore
them

$$\text{if } f(x_0+h) = 0$$

$$\Rightarrow f(x_0) + f'(x_0)h = 0$$

$$\Rightarrow h = \frac{-f(x_0)}{f'(x_0)}$$

if this quadratic
this method terminates in 1 step
and is exact

Newton's method for finding extremum of a function
[find zero of $f'(x)$]

$$x_{n+1} \leftarrow x_n - \frac{f'(x_n)}{f''(x_n)}$$

ex. $f(x) = (x-2)(x-6) = x^2 - 8x + 12$

$$x_0 = 8$$

$$x_1 = 8 - \frac{8}{2} = 8 - 4 = 4 \leftarrow \text{minimum at point. this}$$

consider $F: \mathbb{R}^k \rightarrow \mathbb{R}$

$$F(x_0+h) = F(x_0) + \left(\frac{\partial F}{\partial x} \Big|_{x_0} \right)^T h + \dots$$

ignore higher order terms

$$F: \mathbb{R}^k \rightarrow \mathbb{R}$$

$$J = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \dots & \frac{\partial y_1}{\partial x_k} \\ \vdots & \vdots & & \vdots \\ \frac{\partial y_k}{\partial x_1} & \frac{\partial y_k}{\partial x_2} & & \frac{\partial y_k}{\partial x_k} \end{pmatrix} \quad \left\| \begin{array}{l} J \text{ is different} \\ \text{from change} \\ \text{of p.v.} \end{array} \right.$$

$$h = \begin{pmatrix} h_1 \\ \vdots \\ h_k \end{pmatrix}$$

To find zero of F

$$0 = F(x_0+h) \approx F(x_0)$$

$$+ J|_{x_0} \cdot h$$

$$h = -J|_{x_0}^{-1} F(x_0)$$

To find min of $G: \mathbb{R}^n \rightarrow \mathbb{R}$

and zero of $f = \frac{\partial G}{\partial x}$

In this case $J(G) =$ Matrix of 1st derivatives

$$\text{Hessian} = \begin{pmatrix} \frac{\partial^2 G}{\partial x_1^2} & \frac{\partial^2 G}{\partial x_1 \partial x_2} & \dots & \frac{\partial^2 G}{\partial x_1 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 G}{\partial x_2 \partial x_1} & \dots & \dots & \frac{\partial^2 G}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 G}{\partial x_n \partial x_1} & \dots & \dots & \frac{\partial^2 G}{\partial x_n \partial x_n} \end{pmatrix}$$

$$x_{n+1} \leftarrow x_n - H|_{x_n}^{-1} \frac{\partial G}{\partial x} \Big|_{x_n}$$

$$g(x) = 5x_1^2 + 6x_1x_2 + 5x_2^2$$

find minimum of g

$$\frac{\partial g}{\partial x} = \begin{pmatrix} 10x_1 + 6x_2 \\ 6x_1 + 10x_2 \end{pmatrix}$$

$$H = \begin{pmatrix} 10 & 6 \\ 6 & 10 \end{pmatrix}$$

$$H^{-1} = \frac{1}{24} \begin{pmatrix} 6 & -6 \\ -6 & 10 \end{pmatrix}$$

$$h(x) = 5x_1^3 + 6x_1x_2 + 3x_2^2$$

$$\frac{\partial h}{\partial x} = \begin{pmatrix} 15x_1^2 + 6x_2 \\ 6x_1 + 6x_2 \end{pmatrix}$$

$$H = \begin{pmatrix} 30x_1 & 6 \\ 6 & 6 \end{pmatrix}$$

logistic regression

Also called

iterative reweighted
least squares
(IRLS)

$$\frac{\partial \log L}{\partial w} = \sum_i (b_i - y_i) \phi(x_i)$$

$$\frac{\partial \log L}{\partial w \partial w^T} = - \sum_i \left(\frac{\partial y_i}{\partial w} \right) \phi(x) \frac{\partial y_i}{\partial w^T}$$

$$= - \sum_i \phi(x_i) \frac{\partial y_i}{\partial w^T}$$

$$A_i \begin{pmatrix} \downarrow \\ \sum a_i b_i^T \\ = A^T B^T \end{pmatrix}$$

$$= - \sum_i \phi(x_i) y_i (1 - y_i) \phi(x_i)^T$$

almost
same as calculated
earlier

$$= - \sum_i \phi(x_i) y_i (1 - y_i) \phi(x_i)^T$$

$$= - \Phi^T R \Phi$$

$$R = \text{diag} \{ y_i (1 - y_i) \}_i$$

$R \equiv$ positive (cov)

$$w \leftarrow w + (\Phi^T R \Phi)^{-1} \Phi^T (t y)$$

$$\text{or } w \leftarrow w - (\Phi^T R \Phi)^{-1} \Phi^T (y - t)$$

Proof that our function is concave (i.e. has 1 maximum)

~~or guarantee that~~

we need $H \preceq 0$

negative

positive definite

$$c^T H c \preceq 0$$

c - arbitrary vector

$$-c^T \Phi^T R \Phi c < 0$$

$c \neq 0$

$$-c^T (\Phi^T R^{1/2}) (R^{1/2} \Phi c) < 0$$

$$-\underbrace{(R^{1/2} \Phi^T c)^T (R^{1/2} \Phi c)}_{\text{norm}} < 0$$

this holds if $R^{1/2} \Phi$ is full rank

i.e. if columns in data matrix are linearly independent

Trade off

Newton

optimal step size

but compute hessian

GD

fixed non optimal

step size

but ~~used~~

no additional

computation

Machine Learning

Generative model

K classes $c=1, \dots, c=k$

Prob. of class k : C_k

Prob. of generating data : $P(X|C_k)$

Given model.

prediction

$$\text{Softmax} = \frac{e^{a_k}}{\sum_j e^{a_j}}$$

$$\text{sigmoid}, \sigma = \frac{1}{1 + e^{-(a_1 - a_2)}}$$

$$a_j = \ln \frac{P(C_j)}{P(X|C_j)}$$

When $X|C_j \sim N(\mu_j, \Sigma_j)$

Shared $\Sigma = \Sigma_j \forall j$

Distinct μ_j

Shared Σ , 2 class

$$\text{Prediction } P(C=1) = \sigma(w^T X + w_0)$$

$$\text{let } w = \Sigma^{-1}(\mu_1 - \mu_2)$$

$$w_0 = \ln \frac{P(C_1)}{P(C_2)} - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2$$

use Σ^{-1} is easy

(from data dist. to model)

-ve definite hessian
⇒ concave fn.



Single max

Discriminative model

Logistic Regression

$$P(c=1 | x) = \sigma(w^T \phi(x_i))$$

absorbed bias term w_0 into extra dimension

for analogy $\phi(x) = x$

$$a_i = w^T \phi(x_i)$$

$$y_i = \sigma(a_i)$$

$$t_i \sim \text{Bernoulli}(y_i)$$

Max^m
likelihood
is hard

No closed form solution

⇒ Optimization: Gradient ascent or Newton's method

$$\frac{\partial \log L}{\partial w} = \sum_i \phi(x_i) [t_i - y_i]$$

$$\frac{\partial^2 \log L}{\partial w^2} = -\Phi^T R \Phi$$

$$R = \text{diag} \{y_i (1 - y_i)\}$$

Naive Bayes

Generative model for discrete data.

Say Binary features.

$$\Phi = \begin{pmatrix} 0 & 1 & 1 \\ 0 & 1 & 1 \\ 0 & \dots & \dots \end{pmatrix}$$

$$P(X | c_k) = ?$$

How to write a compact prob. dist.?

- Given Ex
- $P(000)$
 - $P(001)$
 - $P(010)$
 - $P(011)$
 - $P(100)$
 - $P(101)$
 - $P(110)$
 - $P(111)$

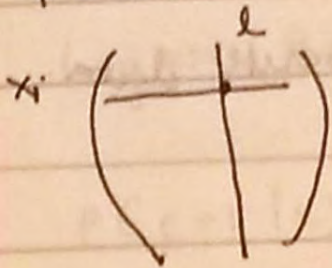
2^n -exponential
for general dist.

Specifying a general dist.
requires too much time / space

A very simple restriction:
assume that features are
conditionally independent
given the label.

$$P(x_i | x_j | c)$$

Notation



x_{il} = l th bit of i th example.

$$P(x_i | C_k) = \prod_{l=1}^d \mu_{kl}^{x_{il}} (1 - \mu_{kl})^{1 - x_{il}}$$

μ_{kl} = parameter for l th bit given that label = k th class

Let first label example have label 3

$$P(C_3) P(x_1 | C_3) = P(C_3) \prod_{l=1}^d \mu_{3l}^{x_{1l}} (1 - \mu_{3l})^{1 - x_{1l}}$$

$$L = \prod_i \prod_k \left[P(C_k) \prod_l \mu_{kl}^{x_{il}} (1 - \mu_{kl})^{1 - x_{il}} \right]$$

boolean $\sum_{i=1}^n$
tick

$$\log L = \sum_i \sum_k \left[\ln P(C_k) + \sum_l x_{il} \ln \mu_{kl} + \sum_l (1 - x_{il}) \ln (1 - \mu_{kl}) \right]$$

tick

MLE for $P(C_k)$ is same as general case

$$\frac{\partial \log L}{\partial \mu_{kl}} = \sum_i \left[\frac{x_{il}}{\mu_{kl}} - \frac{L x_{il}}{1 - \mu_{kl}} \right] = 0$$

tick

$$\Rightarrow \sum_{i=1}^n x_{i1}(1 - \mu_{k1}) - \sum_{i=1}^n \mu_{k1}(1 - x_{i1}) = 0$$

$$\sum_{i=1}^n x_{i1} - \sum_{i=1}^n \mu_{k1} - \sum_{i=1}^n \mu_{k1} + \sum_{i=1}^n \mu_{k1} x_{i1} = 0$$

$$\sum_{i=1}^n x_{i1} - N_k \mu_{k1} = 0$$

$$\sum_{i=1}^n x_{i1} = N_k \mu_{k1}$$

$$\mu_{k1} = \frac{\sum_{i=1}^n x_{i1}}{N_k} = \frac{\text{no. of examples in } P^{\text{th}} \text{ test set}}{\text{no. of examples of class } k}$$

It can be shown that

$$P(c_k | x) = \frac{e^{a_k}}{\sum_j e^{a_j}}$$

$$a_j = w_j^T x + w_{j0}$$

$$\exists w, w_0$$

for some

1st order

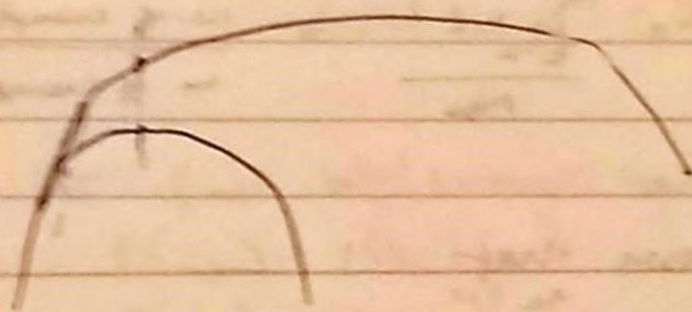
gradient ascent
move with gradient
 $w \leftarrow w + \eta \frac{\partial \log L}{\partial w}$

gradient descent
move against gradient
 $w \leftarrow w - \eta \frac{\partial \log L}{\partial w}$

0th order

Newton's method.

assume that the function is quadratic
and jump to the location that gives the
maximum of the quadratic fn.



2nd order

$$w \leftarrow w + H^{-1} \frac{\partial \log L}{\partial w}$$

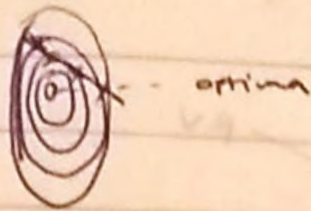
⊙ compute hessian

⊙ invert hessian

$$H = \frac{\partial^2 \log L}{\partial w \partial w^T}$$

alternative: approximate hessian in linear time

Line search. to choose η



want to find ideal η

① Brute force search (expensive)

② backtracking line search

Backtracking Line search
for minimization

initialize $\eta = 1$

$$\text{while } f(w + \eta \frac{\partial f}{\partial w}) > f(w) - \frac{\eta}{2} \left\| \frac{\partial f}{\partial w} \right\|^2$$

$$\eta \leftarrow \frac{\eta}{2}$$

Compromise backtracking line search \hat{w} gradient ascent/descent

What if computing derivative is expensive?

$$\frac{\partial L}{\partial w} = \sum_{i=1}^N \phi(x_i) (f_i - y_i)$$

Sum over

all data points
expensive

Stochastic Gradient Descent

$$GD: w \leftarrow w - \eta G$$

$$G = \frac{\partial f}{\partial w}$$

if w can get \hat{G} ← RV

$$s.t. E[\hat{G}] = G \quad // \text{ unbiased estimator}$$

then we can use \hat{G} instead of G

$\rightarrow w \leftarrow w - \eta_t \hat{G}$ converges to the \min^m of f .

where η_t must satisfy some conditions

$$\eta_t = \frac{1}{t}$$

* Cheap random estimate of gradient for logistic regression

To get \hat{G}

pick $i \in 1 \dots N$ at random (uniformly)

$$\hat{G} = N \phi(x_i) (t_i - y_i)$$

$$E(\hat{G}) = \frac{1}{N} \sum_i N \phi(x_i) (t_i - y_i) = \sum \phi_i (t_i - y_i)$$

Alternative minibatch SGD

take k examples and take mean
instead of 1 example

~~... $(i-1) \theta + \eta \sum_{j=1}^k \nabla_{\theta} \ell_j$...~~

~~... $(i-1) \theta + \eta \sum_{j=1}^k \nabla_{\theta} \ell_j$...~~

~~... $(i-1) \theta + \eta \sum_{j=1}^k \nabla_{\theta} \ell_j$...~~

~~... $(i-1) \theta + \eta \sum_{j=1}^k \nabla_{\theta} \ell_j$...~~

Machine Learning

Exponential family of distributions

Any distribution that can be written in form $p(x) = h(x) g(\eta) e^{\eta^T U(x)}$ is a member of exponential family

* constrain the form of PDF.

$$\eta^T U(x) = A(\eta)$$
$$A(\eta) = -\log g(\eta) \Rightarrow p(x) = h(x) e^{\eta^T U(x) - A(\eta)}$$

Bernoulli \in exp. family

$$p(x) = \mu^x (1-\mu)^{1-x}$$
$$= e^{(\log \mu)x + (\log(1-\mu))(1-x)}$$
$$= e^{(\log \mu)x + \log(1-\mu) - \log(1-\mu)x}$$
$$= \frac{e^{(\log \mu)x}}{e^{\log(1-\mu)x}}$$
$$= (1-\mu) e^{x \log \frac{\mu}{1-\mu}}$$

$h(x)$ = base measure

$g(\eta)$ = log normalizer

η = Natural parameters
canonical parameters

$U(x)$ = Sufficient Statistics

dim. $\eta = 1$ $\eta = \log \frac{\mu}{1-\mu}$ $h(x) = 1$

$g(\eta) = ? = 1 - \mu$

$\mu = \log \frac{\mu}{1-\mu} \Rightarrow e^\eta = \frac{\mu}{1-\mu} \Rightarrow e^\eta - \mu e^\eta - \mu = 0$

$\Rightarrow \frac{e^\eta}{1+e^\eta} = \mu$

$g(\eta) = \frac{e^\eta}{1+e^\eta} = \frac{1}{1+e^{-\eta}}$

$\mu = \frac{1}{1+e^{-\eta}} = \sigma(\eta)$
sigmoid

$P(x) = \frac{1}{1+e^\eta} \cdot e^{\eta x}$
 $\underbrace{\frac{1}{1+e^\eta}}_{g(\eta)}$ $\underbrace{e^{\eta x}}_{U(x)=x}$
 $\eta(x) = 1$

$E[U(x)] = \text{mean parameter} = \theta$ $\eta = \mu$ for bernoulli

if entries of $U(x)$ are linearly independent

then $\theta \xleftrightarrow{1-1} \eta$

Normal

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

$$h(x) = 1 \quad u(x) = \begin{pmatrix} x \\ x^2 \end{pmatrix} \quad \eta = \begin{pmatrix} \eta_1 \\ \eta_2 \end{pmatrix} = \begin{pmatrix} \mu/\sigma^2 \\ -1/2\sigma^2 \end{pmatrix}$$

$$g(\eta) = (2\pi)^{-1/2} \sigma^{-1} \exp\left\{\frac{-\eta^2}{4\eta_2}\right\}$$
$$= \frac{\sqrt{2\eta_2}}{\sqrt{\pi}} \exp\left\{-\eta_1^2/4\eta_2\right\}$$

Poisson

$$p(x) = \frac{\lambda^x e^{-\lambda}}{x!}$$
$$\frac{1}{x!} = \frac{1}{g(\eta) h(x)}$$

$$\frac{1}{x!} \rightarrow \frac{(\log \lambda)^x}{e^x}$$

$$\dim = 1 \quad \eta = \log \lambda \quad \lambda = e^\eta$$

$$h(x) = \frac{1}{x!} \quad u(x) = x$$

Fact 1 If x_1, \dots, x_n are iid sampled from an exp. family distribution then

$$P(x_1, \dots, x_n) = \left(\prod_i h(x_i) \right) (g(\eta))^n e^{\eta^T \sum_{i=1}^n U(x_i)}$$

U : Sufficient stats since only sum is needed for MLE and not each x_i .

Fact 2: $E[U(x)] = -\frac{\partial}{\partial \eta} \log g(\eta)$

$$\text{cov}(U(x)) = E \left[\left(U(x) - E[U(x)] \right) \left(U(x) - E[U(x)] \right)^T \right]$$

$$= -\frac{\partial^2}{\partial \eta \partial \eta^T} \log g(\eta)$$

Fact 3: (S.1) If we have one sample from exp. ^{family} dist. Max likelihood solution is obtained when

$$U(x) = -\frac{\partial}{\partial \eta} \log g(\eta)$$

(S.2) For iid samples

$$\frac{1}{N} \sum_i U(x_i) = -\frac{\partial}{\partial \eta} \log g(\eta)$$

$$\rightarrow \text{Max } L = \frac{1}{N} \sum U(x_i) = E[U(x)]$$

Max L for Bernoulli dist

Use fact 3

$$p(x) = \frac{\lambda^x e^{-\lambda}}{x!}$$

$$h(\eta) = \frac{1}{x!} \eta \log \lambda \quad \lambda = e^\eta \quad U(\eta) = \dots$$

~~max L~~

$$p(x) = \mu^x (1-\mu)^{1-x}$$

$$\eta = \log \frac{\mu}{1-\mu} \quad h(\eta) = 1$$

$$\frac{\partial}{\partial \eta} g(\eta)$$

$$g(\eta) = 1 - \mu = \frac{1 - e^\eta}{1 + e^\eta}$$

$$= \frac{\partial}{\partial \eta} \left(\frac{1}{1 + e^\eta} \right) \\ = - \frac{e^\eta}{(1 + e^\eta)^2}$$

$$= \frac{1}{1 + e^\eta}$$

$$\frac{\partial}{\partial \eta} \log g(\eta)$$

$$\bar{x} = - \frac{1}{g(\eta)} \frac{\partial}{\partial \eta} g(\eta)$$

$$\bar{x} = \frac{e^\eta}{(1 + e^\eta)^2} (1 + e^\eta)$$

$$\bar{x} = \frac{e^\eta}{1 + e^\eta}$$

$$\eta = \log \frac{\bar{x}}{1 - \bar{x}}$$

Sanity check.

$$\mu = \frac{1}{1+e^{-\eta}} = \frac{1}{1+\frac{1-\mu}{\mu}} = \mu$$

Solving using fact 2 + fact 3 together.

$$E[\mu(x)] = \bar{x}$$

$$\Rightarrow \mu = \bar{x}$$

$$E[U(x)] = \frac{1}{N} \sum U(x_i)$$

$$\begin{pmatrix} E[x] \\ E[x^2] \end{pmatrix} = \frac{1}{N} \begin{pmatrix} \sum x_i \\ \sum x_i^2 \end{pmatrix}$$

$$\begin{pmatrix} E[x] \\ E[x^2] \end{pmatrix} = \frac{1}{N} \begin{pmatrix} \sum x_i \\ \sum x_i^2 \end{pmatrix}$$

$$\mu = \frac{1}{N} \sum x_i$$

$$\mu^2 + \sigma^2 = \frac{1}{N} \sum x_i^2$$

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

Proof of Fact 3.1

$$L = h(x) g(\eta) e^{\eta^T u(x)}$$

$$\log L = \log h(x) + \log g(\eta) + \eta^T u(x)$$

$$\frac{\partial \log L}{\partial \eta} = \frac{\partial}{\partial \eta} \log g(\eta) + u(x) = 0$$

$$u(x) = -\frac{\partial}{\partial \eta} \log g(\eta)$$

Fact 3.2

$$L = \prod_i h(x_i) g(\eta)^N e^{\eta^T \sum u(x_i)}$$

$$\log L = \sum \log h(x_i) + N \log g(\eta) + \eta^T \sum u(x_i)$$

$$\frac{\partial \log L}{\partial \eta} = \left(N \frac{\partial}{\partial \eta} \log g(\eta) + \sum u(x_i) \right) = 0$$

Proof of fact 2

$$\int h(x) g(\eta) e^{\eta^T u(x)} dx = 1$$

$$g(\eta) \int h(x) e^{\eta^T u(x)} dx = 1$$

$$\frac{\partial}{\partial \eta} (g(\eta)) \left(\int h(x) e^{\eta^T u(x)} dx \right) = 0$$

product rule

$$\frac{\partial}{\partial \eta} g(\eta) \int h(x) e^{\eta^T u(x)} dx + g(\eta) \frac{\partial}{\partial \eta} \left(\int h(x) e^{\eta^T u(x)} dx \right) = 0$$

$$\int h(x) e^{\eta^T U(x)} dx = \frac{1}{g(\eta)} \quad \text{and} \quad \int h(x) e^{\eta^T U(x)} U(x) dx = \frac{E[U(x)]}{g(\eta)}$$

$$\frac{\partial g(\eta)}{\partial \eta} \frac{1}{g(\eta)} + \cancel{g(\eta)} \frac{E[U(x)]}{\cancel{g(\eta)}} = 0$$

$$\frac{1}{g(\eta)} \frac{\partial g(\eta)}{\partial \eta} + E[U(x)] = 0$$

$$\frac{\partial \log g(\eta)}{\partial \eta} + E[U(x)] = 0 \quad \Rightarrow \quad E[U(x)] = -\frac{\partial \log g(\eta)}{\partial \eta}$$

$$\text{cov}(U(x_i), U(x_j)) = E[(U(x_i) - E[U(x_i)])(U(x_j) - E[U(x_j)])]$$

$$\frac{\partial \log g(\eta)}{\partial \eta} + g(\eta) \int h(x) e^{\eta^T U(x)} U(x) dx = 0$$

$$\frac{\partial \log g(\eta)}{\partial \eta_i} + g(\eta) \int h(x) e^{\eta^T U(x)} U_i(x) dx = 0$$

$$\frac{\partial^2 \log g(\eta)}{\partial \eta_i \partial \eta_j} + g(\eta) \int h(x) e^{\eta^T U(x)} U_i(x) U_j(x) dx = 0$$

$$+ \frac{\partial g(\eta)}{\partial \eta_j} \int h(x) e^{\eta^T U(x)} U_i(x) dx = 0$$

$$\frac{\partial}{\partial \eta_k} (\log g(\eta)) + E(U_k(x)) = 0$$

$$E(U_k(x)) = - \frac{\partial}{\partial \eta_k} (\log g(\eta))$$

$$= - \frac{g(\eta)}{g(\eta)} + \frac{1}{g(\eta)} \frac{\partial g(\eta)}{\partial \eta_k}$$

$$\frac{\partial}{\partial \eta_k} g(\eta) = E(U_k(x)) g(\eta)$$

$$\int h(x) e^{\eta^T U(x)} U_i(x) dx = \frac{E[U_i(x)]}{g(\eta)}$$

$$\int h(x) e^{\eta^T U(x)} U_i(x) U_k(x) dx = \frac{E[U_i(x) U_k(x)]}{g(\eta)}$$

$$\frac{\partial^2 \log g(\eta)}{\partial \eta_i \partial \eta_k} = E[U_i(x)] E[U_k(x)] + E[U_i(x) U_k(x)] = 0$$

$$- E[U_i(x)] E[U_k(x)] + E[U_i(x) U_k(x)] = 0$$

$$= \frac{\partial^2}{\partial \eta_i \partial \eta_k} [\log g(\eta)]$$

$$L \propto g(\eta)^N e^{\eta^T \sum_i U(x_i)}$$

no. of obs- \swarrow

\searrow sum of sufficient stats

Conjugate prior should have

$$p(\eta) \propto g(\eta)^{\nu} e^{\eta^T g}$$

$$= g(\eta)^{\nu} e^{g^T \eta} = g(\eta)^{\nu} e^{(\nu \bar{g})^T \eta}$$

$\nu \equiv$ pretend
to have seen
 ν no. of
examples

$\bar{g} =$ mean
of such
pseudo
observations

added
counts

mean
for newer
obs.

$$p(x|y) = \frac{1}{Z} \sum_{\theta} p(\theta) p(x|y, \theta)$$

Machine Learning

Quiz 3 model selection - GLM, HW3.

Generalised Linear Models

Linear Regression

$$a_i = w^T \phi(x_i)$$

$$t_i \sim N(a_i, \frac{1}{\beta})$$

Logistic Regression

$$a_i = w^T \phi(x_i)$$

$$y_i = \sigma(a_i)$$

$$t_i \sim \text{Bernoulli}(y_i)$$

Poisson Regression

$$a_i = w^T \phi(x_i)$$

$$y_i = e^{a_i}$$

$$t_i \sim \text{Poisson}(y_i)$$

$$\frac{\partial \log L}{\partial w} = \sum (t_i - y_i) \phi(x_i)$$

$$\frac{\partial^2 \log L}{\partial w^2} = - \sum y_i (1 - y_i) \phi(x_i) \phi(x_i)^T$$

$$R = - \Phi^T R \Phi$$

$$R = \text{diag}(y_i (1 - y_i))$$

Poisson Regression

$$L = \prod_i y_i^{t_i} e^{-y_i} \frac{1}{t_i!}$$

$$\log L = \sum t_i \log y_i - y_i - \log t_i!$$

$$= \sum t_i (\log e^{a_i}) - y_i - \log t_i!$$

$$= \sum t_i a_i - y_i - \log t_i!$$

$$\frac{\partial \log L}{\partial w} = \sum t_i \phi(x_i) - e^{a_i} \phi(x_i)$$

$$= \sum (t_i - y_i) \phi(x_i)$$

same form.

$$\begin{aligned}
 \frac{\partial \log L}{\partial w^T} &= \sum_i \phi(x_i) \frac{\partial y_i}{\partial w^T} \\
 &= \sum_i \phi(x_i) e^{y_i} \phi(x_i)^T \\
 &= \sum_i e^{y_i} \phi(x_i) \phi(x_i)^T \\
 &= \sum_i y_i \phi(x_i) \phi(x_i)^T \\
 &= \Phi^T R \Phi \quad R = \text{diag}(y_i)
 \end{aligned}$$

$$w \leftarrow w - (\Phi^T R \Phi)^{-1} \Phi^T (y - t)$$

Exponential Family of Distributions

$$P(x|\eta) = h(x) g(\eta) e^{\eta^T U(x)}$$

η : natural parameter

g : normalizer

U : sufficient statistics

$\mu = E[U(x)]$: mean parameter

If U is linearly independent then

we can write PDF in two ways (given η with μ)

$$\eta = \Psi(\mu)$$

$$\mu = \Psi^{-1}(\eta)$$

Bernoulli

$$P(x|\eta) = \eta^x (1-\eta)^{1-x}$$

$$P(x|\eta) = \frac{1}{1+e^{-\eta}} e^{\eta x}$$

$$\eta = \frac{1}{1+e^{-\eta}} = \sigma(\eta)$$

$$\eta = \log \frac{\eta}{1-\eta}$$

Poisson

$$P(x|\eta) = \frac{1}{x!} e^{-e^{\eta}} e^{\eta x}$$

$$\eta = \log \lambda$$

$$\lambda = e^{\eta}$$

Any PDF $f(x)$ can be changed to add a scale parameter.

$$\textcircled{1} P(x) = \frac{1}{s} f\left(\frac{x}{s}\right)$$


Apply to an ID exp. family dist. $U(x) = x$

$$P(x) = \frac{1}{s} h\left(\frac{x}{s}\right) g(\eta) e^{\frac{1}{s} \eta \cdot x}$$

For any exp family dist.

$$E[U(x)] = -s \frac{\partial}{\partial \eta} \log g(\eta)$$

$$\textcircled{2} E[x] = -s \frac{\partial}{\partial \eta} \log g(\eta)$$

scale (wavy line) \rightarrow 

$$\text{cov}(\eta) = -\frac{\partial^2}{\partial \eta \partial \eta^T} \log g(\eta)$$

$$\text{var}(x) = -s^2 \frac{\partial^2}{\partial \eta^2} \log g(\eta)$$

$$p(x) = (2\pi)^{-1/2} \sigma^{-1/2} \exp\left\{-\frac{1}{2\sigma^2}(x^2 - 2\mu x + \mu^2)\right\}$$

$$\underbrace{(2\pi\sigma^2)^{-1/2} e^{-\frac{1}{2\sigma^2}\mu^2}}_{g(\eta)} \underbrace{e^{-\frac{1}{2\sigma^2}x^2}}_{h(x)} \underbrace{e^{\frac{1}{\sigma^2}\mu x}}_{e^{\frac{1}{\sigma^2}\eta x}}$$

In this representation $\eta = \mu$ $\psi = \text{identity}$

Generalized Linear Model

$$a_i = w^T \phi(x_i)$$

$y_i = f(a_i)$ \leftarrow mean parameter and mean of t_i

f : activation $\eta_i = \psi(y_i)$ natural parameter

f^{-1} : link function.

$$t_i \sim P(t_i | \eta_i) \quad (\text{exp. family dist.})$$

1D exp family dist.

Canonical link function picks $f(a) = \psi^{-1}(a)$

$$\eta_i = \sum \psi(y_i) = \psi(\psi^{-1}(a_i)) = a_i$$

for linear regression, $\psi = \text{identity}$.

for logistic regression

$$t_i \sim \text{Bernoulli}(y_i) \\ \sim \text{Ber.}^{\text{natural}}(t_i | a_i)$$

for poisson

$$t_i \sim \text{Poisson}(y_i) \\ \sim \text{natural poisson}(t_i | a_i)$$

GLM

$$L = \prod_i \frac{1}{s} h\left(\frac{t_i}{s}\right) g(\eta_i) e^{\frac{1}{s} \eta_i t_i}$$

$u(x) = x$

$t_i \sim \text{ID exp.}$

family dist.

to scale

$$\log L = \sum_i \log s + \log h\left(\frac{t_i}{s}\right) + \log g(\eta_i) + \frac{1}{s} \eta_i t_i$$

$$\frac{\partial \log L}{\partial w} = \sum_i \frac{1}{g(\eta_i)} \frac{\partial g(\eta_i)}{\partial \eta} + \frac{1}{s} \frac{\partial \eta_i}{\partial w} t_i$$

$$= \sum_i \frac{1}{g(\eta_i)} \frac{\partial g(\eta_i)}{\partial \eta} \cdot \frac{\partial \eta_i}{\partial y} \frac{\partial y_i}{\partial a_i} \frac{\partial a_i}{\partial w}$$

$$+ \frac{1}{s} \frac{\partial \eta_i}{\partial y} \frac{\partial y_i}{\partial a_i} \frac{\partial a_i}{\partial w}$$

$$= \sum_i \frac{\partial \log g(\eta_i)}{\partial \eta_i} \phi'(y_i) f'(a_i) \phi(x_i) + \frac{1}{s} t_i \psi(y_i) f'(a_i) \phi(x_i)$$

Canonical link

$$\eta_i = a_i$$

$$\frac{\partial \eta_i}{\partial w} = \phi(x_i)$$

$$\Rightarrow \frac{\partial \log L}{\partial w} = \frac{\partial \log g(\eta_i) \phi(x_i) + \frac{1}{S} t_i \phi(x_i)}{\partial \eta_i}$$

$$E[U(x)] = - \frac{\partial \log g(\eta)}{\partial \eta}$$

$$E[x] = -S \frac{\partial \log g(\eta)}{\partial \eta}$$

$$\frac{\partial \log g(\eta)}{\partial \eta} = -\frac{1}{S} E[x]$$

$$\Rightarrow \frac{\partial \log L}{\partial w} = -\frac{1}{S} E[t] \phi(x_i) + \frac{1}{S} t_i \phi(x_i)$$

$$= -\frac{1}{S} y_i \phi(x_i) + \frac{1}{S} t_i \phi(x_i)$$

$$= \frac{1}{S} \sum_i (t_i - y_i) \phi(x_i)$$

choose f

ψ is determined by exp-family dist. type

t_i describes data-generation

$$\frac{\partial \log L}{\partial w} = -\frac{1}{S} \sum f'(a_i) \phi(x_i) \phi(x_i)^T$$

$$= -\frac{1}{S} \phi^T R \phi$$

$$R = \text{diag}(r_i) \quad r_i = f'(a_i)$$

$$\frac{\partial \log L}{\partial w} = \sum_i \frac{\partial \log g(y_i)}{\partial \eta_i} \psi'(y_i) f'(a_i) \phi(x_i)$$

$$= \sum_i \frac{1}{s} (t_i - y_i) \psi'(y_i) f'(a_i) \phi(x_i)$$

ϕ^T

$$\frac{\partial \log L}{\partial w \partial w^T} = \sum_i \frac{1}{s} (t_i - y_i) [\psi'(y_i) f''(a_i) \phi(x_i) \phi(x_i)^T$$

$$+ f'(a_i)^2 \psi''(y_i) \phi(x_i) \phi(x_i)^T]$$

$$- \frac{1}{s} \psi'(y_i) f'(a_i)^2 \phi(x_i) \phi(x_i)^T$$

$$\text{let } r_i = (t_i - y_i) [\psi'(y_i) f''(a_i) + f'(a_i)^2 \psi''(y_i) - \psi'(y_i) f'(a_i)^2]$$

$$\Rightarrow \sum_i \frac{1}{s} r_i \phi(x_i) \phi(x_i)^T = -\frac{1}{s} \Phi^T R \Phi$$

$$w \leftarrow w + s (\Phi^T R \Phi)^{-1} \frac{1}{s} \Phi^T (t - y)$$

} Canonical
link

$$R = \text{diag}(f'(a_i))$$

Normal $f(a) = a$ $f'(a) = 1$ $R = I$

Bernoulli $f(a) = \sigma(a)$ $f'(a) = y_i(1 - y_i)$

Poisson $f(a) = e^a$ $f'(a) = e^a$

Machine Learning.

Logistic Regression

$$a_i = w^T \phi(x_i)$$

$$y_i = \sigma(a_i)$$

$$L = \prod_i \sigma(w^T \phi(x_i))^{t_i} (1 - \sigma(w^T \phi(x_i)))^{1-t_i}$$

$$\frac{\partial \log L}{\partial w} = \Phi^T (t - y)$$

$$\frac{\partial^2 \log L}{\partial w \partial w^T} = -\Phi^T R \Phi \quad R = \text{diag}(y_i (1 - y_i))$$

$$w \leftarrow w - (\Phi^T R \Phi)^{-1} \Phi^T (t - y)$$

iterative method

what would be a conjugate prior for w ?

prior \times likelihood \propto posterior

same form

No known PDF satisfies conjugate requirement

Let's use a normal prior, $w \sim N(\mu_0, \Sigma_0)$

(since it's most convenient)



Cheapest approx. by gaussian

This model is not conjugate

$$t_i | w \sim \text{Bernoulli}(w^T \phi(x_i), \frac{1}{\beta})$$

Posterior \propto prior \times likelihood

$$\text{also, posterior} \propto \prod y_i^{t_i} (1-y_i)^{1-t_i} \left(\frac{1}{2\pi}\right)^{d/2} \frac{1}{|S_0|^{d/2}} e^{-\frac{1}{2}(w-w_0)^T S_0^{-1} (w-w_0)}$$

No known form for posterior

Sol 1: Maybe we can represent it "indirectly"
(typically via sampling)

Sol 2: Approximate posterior "as well as we can"

Before Posterior, let's compute MAP.

$$\log \text{Posterior} = \text{const} + \sum t_i \ln y_i + \sum (1-t_i) \ln(1-y_i) - \frac{1}{2} (w-w_0)^T S_0^{-1} (w-w_0)$$

diff wrt w .

$$\frac{d}{dw} \log \text{Posterior} = 0 + \phi^T (t-y) - S_0^{-1} (w-w_0)$$

(no closed form solution for w)

$$\frac{d^2}{dw dw^T} \log \text{Posterior} = -\phi^T R \phi - S_0^{-1}$$

$$\frac{d^2}{dw dw^T} \left(\frac{1}{N} \sum \phi \phi^T \right) = \frac{1}{N} \sum \phi \phi^T$$

$$w \leftarrow w - H^{-1} G$$

$$= w - (H^{-1}) G$$

$$= w - (\phi^T R \phi + S_0^{-1})^{-1} (\phi^T (y-t) + S_0^{-1} (w-w_0))$$

finds w_{MAP}

Laplace Approximation

1D case:

want to approximate $f(x)$ as \mathcal{N}

$$f(x) \propto e^{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}}$$

So, approx. $g(x) = \log f(x)$ as a quadratic fn.

$$x = x_0 + h$$

Taylor expansion \rightarrow

$$g(x) \approx g(x_0) + \underbrace{g'(x_0) \cdot h}_0 + \frac{1}{2} g''(x_0) h^2 + \dots$$

ignore higher order terms

$$= g(x_0) + \frac{1}{2} g''(x_0) (x-x_0)^2$$

$$g'(x_0) = 0 \Rightarrow x_0 = \mu$$

$$A = \frac{1}{\sigma^2} = -g''(x_0)$$

point of i.e. mean of normal = point of curvature at mean (maxima of g is same)

$$f(x) \propto e^{-\frac{A}{2} (x-x_0)^2} = e^{-\frac{1}{2\sigma^2} (x-\mu)^2}$$

Multivariate case

$$A = -H^{-1}$$

$$f(x) \propto e^{-\frac{1}{2} (x-x_0)^T H (x-x_0)}$$

$$= e^{-\frac{1}{2} (x-x_0)^T A^{-1} (x-x_0)}$$

A is known as covariance matrix

Coming back to Logistic Regression

$$x_0 = w_{\text{MAP}}$$

$$H = -(\Phi^T R \Phi + S_0^{-1})$$

$$\Rightarrow m_N = w_{\text{MAP}}$$

$$S_N = [\Phi^T R \Phi + S_0^{-1}]^{-1}$$

$$w_{\text{MAP}} \sim N(m_N, S_N)$$

How to predict?

Using w_{MAP} : Given x_{N+1} predict $p(c=1) = \sigma(w_{\text{MAP}}^T \Phi(x_{N+1}))$

For Bayesian Solution

$$p(c=1) = \int_w p(w | m_N, S_N) p(c=1 | w) dw$$

$$\int_w N(w | m_N, S_N) \sigma(w^T \Phi(x_{N+1})) dw$$

Approximate $\sigma(a)$ as $\phi(\lambda a)$ $\phi = \text{PDF of normal dist}$
mean = 0, variance = 1

$$\lambda = \sqrt{\frac{\pi}{8}}$$

$$\sigma(a) \approx \phi(\lambda a)$$

$$a = w^T \Phi(x_{n+1})$$

$$a \sim N(\Phi(x_{n+1})^T \mu_N, \Phi(x_{n+1})^T \Sigma_N \Phi(x_{n+1}))$$

if w is gaussian, a linear transform of w is also gaussian

$$P(c=1) = \int N(a | \mu_a, \sigma_a^2) \sigma(a) da$$
$$= \int N(a | \mu_a, \sigma_a^2) \Phi(\lambda a) da$$

Final result

$$\sigma \left(\frac{\mu_a}{\sqrt{1 + \frac{\pi}{8} \sigma_a^2}} \right)$$

$$\sigma: \quad x > 0 \Rightarrow 1$$
$$x < 0 \Rightarrow 0$$

(+ same cost of +ve & -ve labels)

$P(c=1)$ depends only on sign of μ_a

Machine Learning

Applying ML Algorithms,

Nearest neighbours and decision trees,

KNN algorithm for classification

store all examples

find k nearest neighbours

predict label based on voting for each neighbour

can be for regression

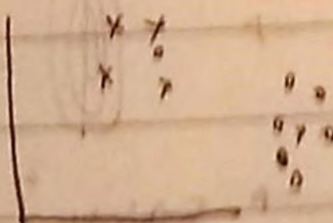
take mean of k nearest neighbours.

* NON PARAMETRIC METHOD.

no prior commitment to hypothesis

Complexity increases as n increases

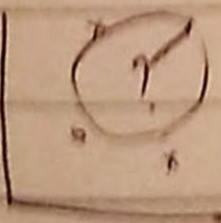
! noise



increase value of k for random noise

! query by example

Search is expensive



use distance inequality

geometric D.P. (k nearest)

k is free parameter

how to choose k ?

use validation data to evaluate how each k performs
model selection

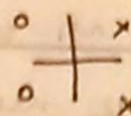
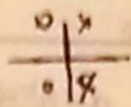
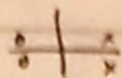
! sensitivity to how data is presented (scale)

since it depends on distance

→ normalize

linear

Z scaling



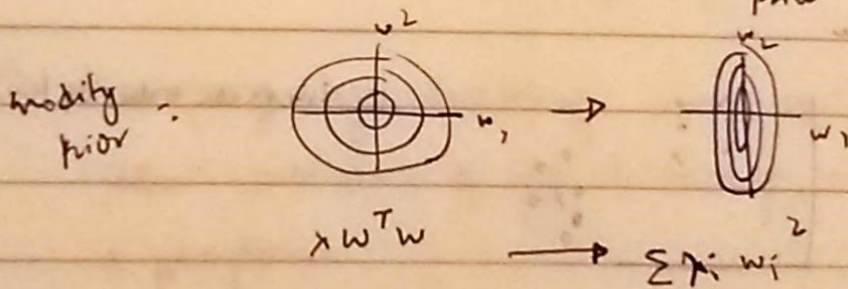
Sensitive to irrelevant features

irrelevant features dominate relevant features

→ apply dimensionality reduction

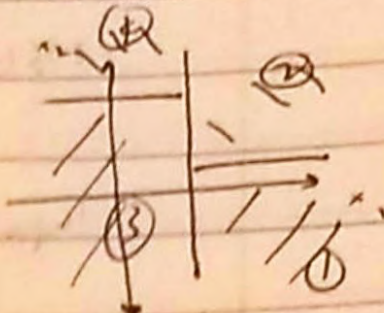
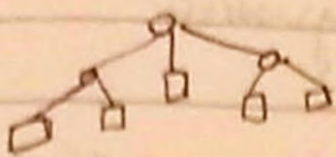
Bayes-LP

$$\sum (t_i - u^T \phi(x_i))^2 + \underbrace{\lambda W^T W}_{\text{prior}}$$



Decision Trees.

non-parametric and non-probabilistic classifier.



recursively split feature space



! tree can be very large (exponential)
! it may overfit

how to build a good tree (best tree = NP hard)

Splitting criteria.

algorithm:

if data has pure class

make leaf

else:

pick feature to split on

split data into subsets

apply algorithm recursively for each subset

Measure uncertainty

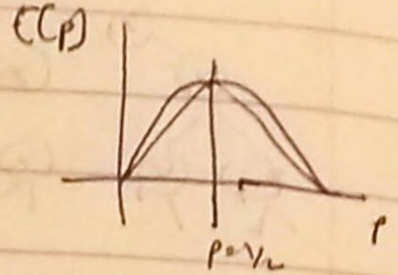
entropy - classification

MSE - regression

(accuracy does not work as well)

Information gain = reduction in uncertainty due to split.

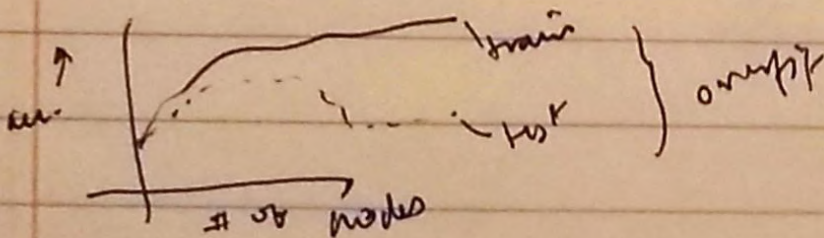
$$\text{Entropy}(P_1, \dots, P_n) = -\sum p_i \log p_i$$



$$\text{Gain}(\text{split}) = \text{Ent}(S) - \sum \frac{|S_j|}{|S|} \text{Ent}(S_j)$$

Real valued Attributes.

x_1									} expanded overfit
2.3	T	2.3	3.5	1.6	2.7	8.3	...		
3.5	-								
1.6	-	-	+	+	-	+	-		
2.7	+	1.6	2.3	2.7	3.5	6.5	8.3		
8.3	-	① only finite # points make d'Hermite							
6.5	+	② no need to test values in [2.3, 2.7]							



solutions to overfitting
prevent -

min # points at any level

min. information gain

grow tree to full size.

then prune use valid. set

} best

ex. Reduced error pruning

App Machine Learning

Applying ML

individual feature preprocessing

linear scaling to $[0, 1]$

$$x \leftarrow \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

normal scaling

$$x \leftarrow \frac{x - \mu}{\sigma}$$

Discretizing features

unsupervised

1. equal bin size predetermined by range
2. equal frequency adapts to data
3. cluster

Supervised.

use labels. Run decision tree on single feature
most helpful after pruning

discrete to numerical

unit vector 0000, 0100, 0010, 0001

increasing weight vectors

1000, 1100, 1110, 1111

Manifold methods

data resides on a "manifold"

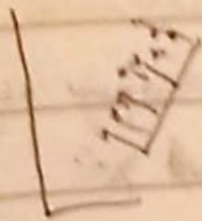
embed data in low dim space, preserving local distances.

process "embedding"

PCA - Principal Component Analysis

linear dim. reduction

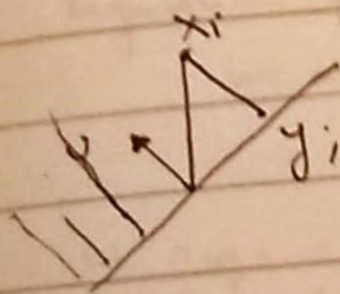
project data onto k dim. with max variance



center data matrix

$$\Phi^T \Phi = V \Lambda V^T \quad \text{eigen decomposition}$$

take top k eigenvectors.



$$\text{Data} = x_1 \dots x_N$$

$$y = u^T x_i$$

$$\bar{x} = \frac{1}{N} \sum x_i$$

$$\bar{y} = \frac{1}{N} \sum y_i$$

$$\bar{y} = \frac{1}{N} \sum u^T x_i = \frac{1}{N} u^T \bar{x}$$

Projected variance

$$J = \frac{1}{N} \sum (y_i - \bar{y})^2$$

want to maximize

vector

scalar

$$\begin{aligned}
 \text{also, } \bar{S} &= \frac{1}{N} \sum_i [U^T (x_i - \bar{x})]^2 \\
 &= \frac{1}{N} \sum_i U^T (x_i - \bar{x}) (x_i - \bar{x})^T U \\
 &= U^T \left[\frac{1}{N} \sum_i (x_i - \bar{x}) (x_i - \bar{x})^T \right] U \\
 &= U^T \Sigma U
 \end{aligned}$$

constraint norm $U = 1$

Use Lagrange Multiplier

objective + λ (constraint) = new objective
 solve for original vars. and λ

$$\max U^T S U \quad \text{s.t. } U^T U = 1$$

$$\equiv U^T S U + \lambda (U^T U - 1)$$

$$\frac{\partial \mathcal{L}(U, \lambda)}{\partial \lambda} = 1 - U^T U = 0 \Rightarrow \|U\| = 1$$

$$\frac{\partial \mathcal{L}(U, \lambda)}{\partial U} = 2SU - 2\lambda U = 0$$

$$SU = \lambda U \quad \begin{array}{l} \text{eigenvalue} \\ \text{eigenvector} \end{array} \quad \|U\| = 1$$

which eigenvector?

$$\int_0^{\infty} U^T S U = U^T \lambda U = \lambda \|U\|^2$$

pick max eigen value

Projection = ΦU

Feature Selection.

Filter method

calculate score (ex. info gain, correlation w/ label)

for each feature

we pick top k

Issue - may end up choosing same / similar duplicate features

L1 Regularization

Regularized linear regression

$$w = \arg \min \underbrace{\frac{1}{2} \sum (\omega^T \phi(x_i) - t_i)^2}_{\text{loss}} + \underbrace{\lambda \sum_{k=1}^K \omega_k^2}_{\text{regularization}}$$

L1 reg.

$$w = \arg \min \frac{1}{2} \text{LOSS} + \lambda \sum_{k=1}^K |w_k|$$

use Laplace distribution

instead of gaussian for prior

Evaluating ML outcomes.

what to measure?

Regression, Classification

MSE

acc.

Confusion matrix

	+	-	← classified as
+	TP	FN	
-	FP	TN	

$$acc = \frac{TP + TN}{TP + TN + FN + FP}$$

IR terminology

Ignore
true
negative

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

of those that I predicted,
how many are POSITIVE

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

of the ones that I should've
found how many did I find

$$F = \frac{2 \cdot PR}{R + P}$$

medical terminology

$$sensitivity = recall$$

accuracy in + class

$$Specificity = \frac{TN}{TN + FP}$$

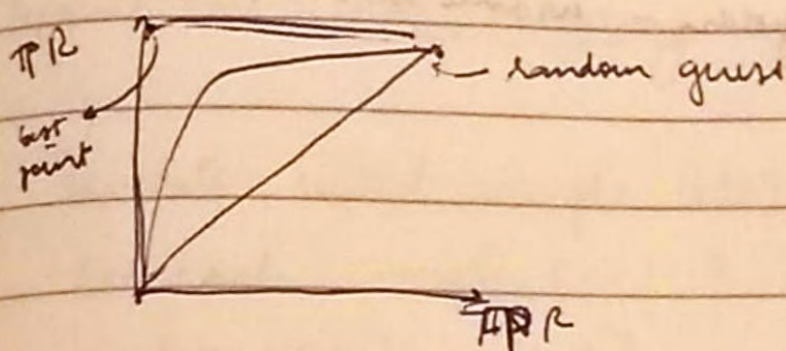
accuracy in - class

signal detection

$$TP_{rate} = Recall$$

$$FP_{rate} = 1 - Specificity$$

ROC (receiver operator characteristic) curve



get points by changing threshold

area under ROC curve

($< 1 \sim$ probability)

\neq that random example from test set is classified correctly

How to measure?

validation set method

+ unbiased estimate of data

- variance due to choice of valid set

- wastes data

do many times and take average

(-introduces bias)

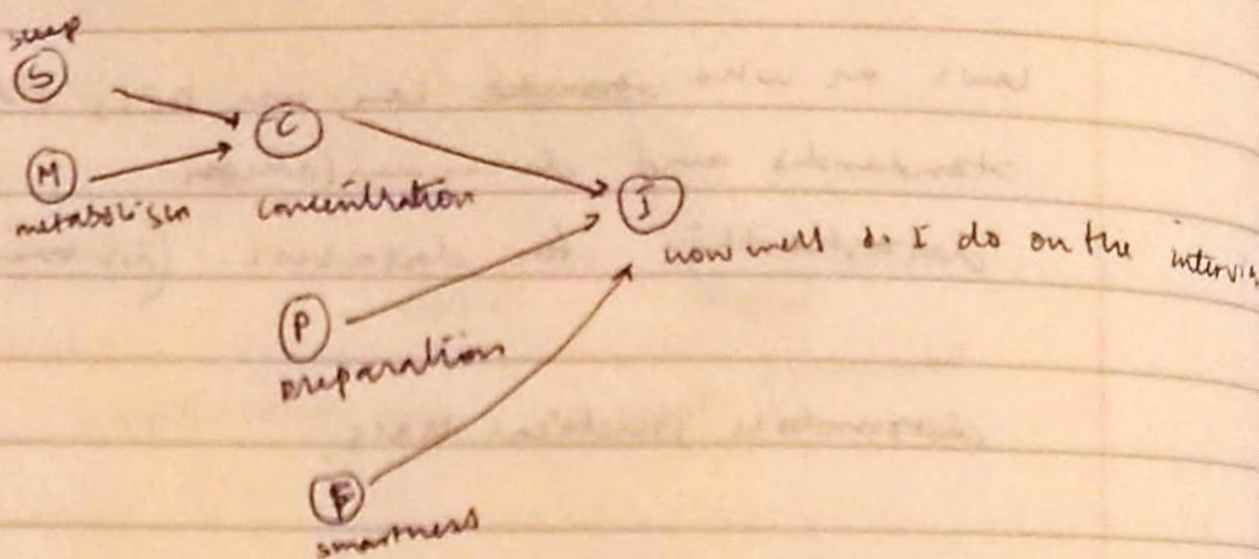
\rightarrow k fold cross validation (disjoint)

(-class label dist. bias)

\rightarrow stratified CV

Machine Learning

Graphical Models



Bayesian network

a directed acyclic graph
probabilistic relationship b/w variables
for every node v : $P(v | \text{parents}(v))$

General form

Nodes are x_1, \dots, x_N

$$P(x_i | \text{Pa}(x_i))$$

Joint Dist

$$P(x_1, \dots, x_N) = \prod_i P(x_i | \text{Pa}(x_i))$$

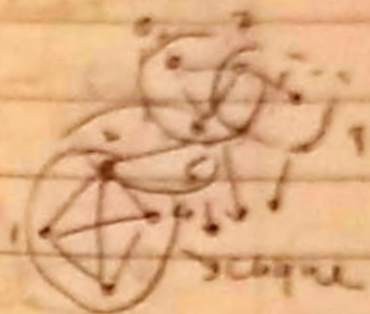
Undirected Graphical models

(Markov Random Fields)

Graph with no edges

For every clique in graph

potential function $\psi(x_c)$



$$P(x_1, \dots, x_n) \propto \prod_c \psi_c(x_c)$$

$$c_1 = 1, 2, 3, 4$$

$$c_2 = 2, 5$$

$$c_3 = 3, 4, 5$$

$$c_4 = 2, 3, 4$$

or



$$c_1 = x_1, x_2$$

$$c_2 = x_2, x_3$$

	x_1	x_2	ψ_c
ψ_{c_1}	0	0	5
	0	1	1
	1	0	3
	1	1	100
	x_2	x_3	
ψ_{c_2}	0	0	1
	0	1	1
	1	0	2
	1	1	2

$$P(x_1, \dots, x_n) = \frac{1}{Z} \prod_c \psi_c(x_c)$$

$$Z = \sum_{x_1, \dots, x_n} \prod_c \psi_c(x_c)$$

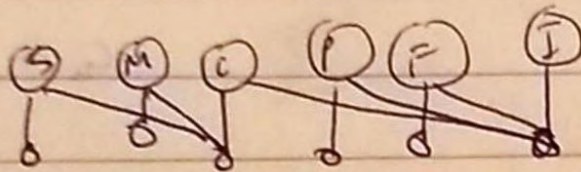
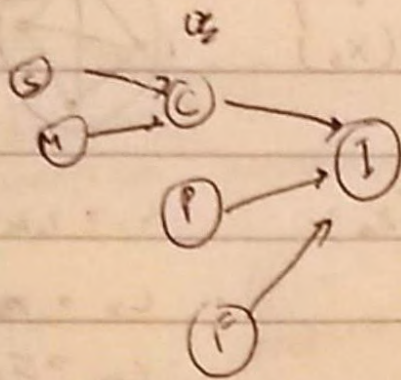
a generalization
of binary net
 $c_i = i$ and parents

$$\psi(x_i, p_i(x_i))$$

More general \rightarrow factor graph.

Random Variables (x_1, \dots, x_n)

and functions of R.V.s.



in linear regression

in Bayesian linear regression

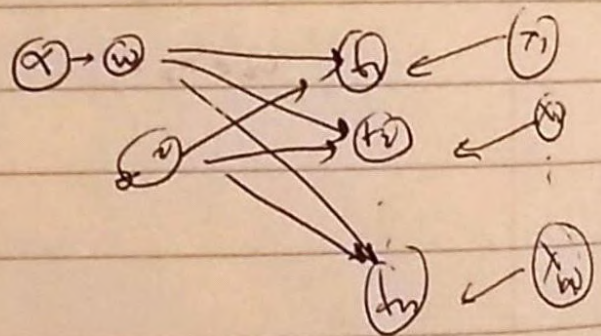
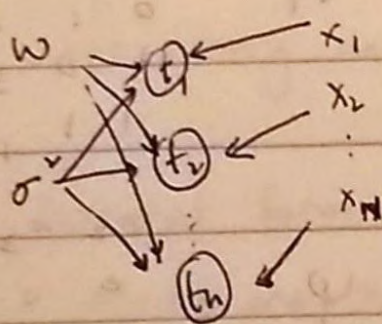
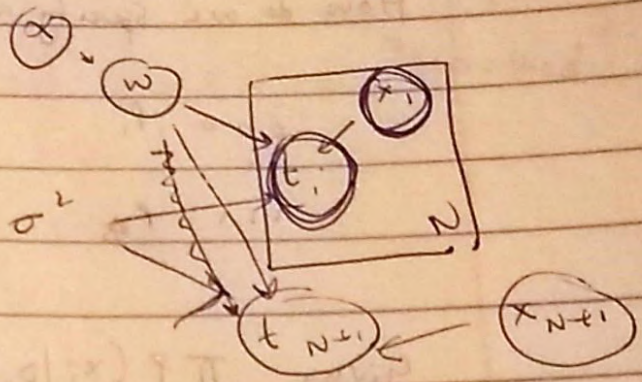
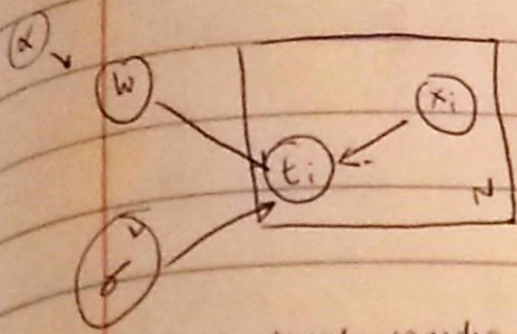


Plate notation.



observe some nodes,
ML, MAP:

find assignment to some variable (w)
s.t. $P(\text{evidence} | w)$ is Max.

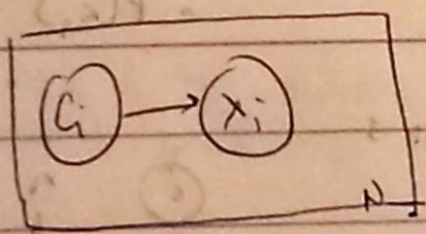
PA

Predictive distribution

find $P(t_{N+1} | \text{Evidence})$

Generative model:

$$P(c_i) \cdot P(x_i | c_i)$$



How do we specify a dist. over 3 binary r.v.?

$$\begin{matrix} 000 & p_1 \\ \vdots & \\ \dots & p_8 \end{matrix} \quad \sum p_i = 1$$

Given $\prod P(x_i | p_i(x_i))$

Q obs ~~$x_3 = 1$~~ $x_3 = 1$

$$P(x_2 = 1) = ?$$

$$P(x_2 = 1) = \sum_{x_1, x_3} P(x_1, x_2 = 1, x_3)$$

Can ~~we~~ first write a full table of joint dist. and marginalize variables not of interest (x_1, x_3)

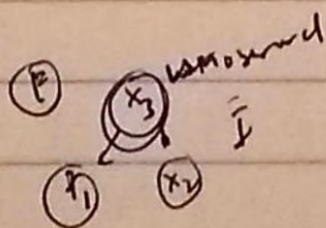
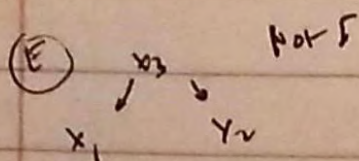
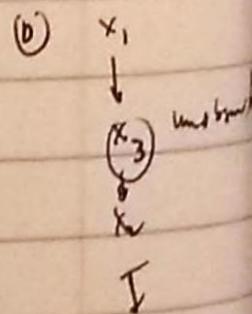
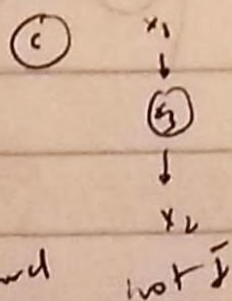
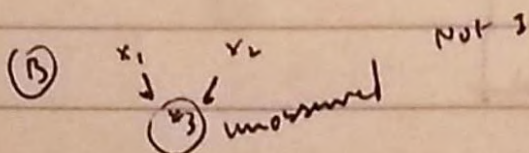
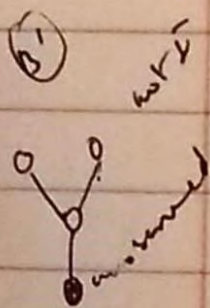
$J = x_1, x_2$

(A) x_1, x_2, x_3 \downarrow x_3

$$P(x_1, x_2, x_3) = P(x_1) P(x_2) P(x_3 | x_1, x_2)$$

$$P(x_1, x_2) = \sum P(x_1) P(x_2) P(x_3 | x_1, x_2)$$

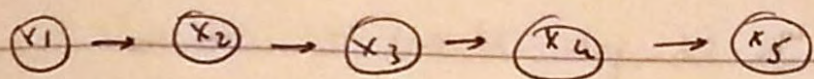
$$= P(x_1) P(x_2) \sum \underbrace{P(x_3 | x_1, x_2)}_{=1}$$



head to head assumed
 incoming outgoing unassumed
 tail to tail unassumed

Theorem: v is independent of U

\Leftrightarrow every path from $v \in V$ to $u \in U$ is "blocked"



$$P(x_1 = 1) = 0.7$$

$$P(x_i = 1 \mid x_{i-1} = 1) = 0.9$$

$$P(x_i = 1 \mid x_{i-1} = 0) = 0.3$$

① $P(x_3 = 1)$

no witness.

② $P(x_3 = 1 \mid x_1 = 1)$

parent to children

③ $P(x_3 = 1 \mid x_5 = 1)$

children to parent

$$P(x_3) = \sum_{x_1, x_2, x_4, x_5} P(x_1) P(x_2 \mid x_1) P(x_3 = 1 \mid x_2) P(x_4 \mid x_3 = 1) P(x_5 \mid x_4)$$

$$P(x_1 = 1) = 0.7$$

$$P(x_2 = 1) = \sum_{v \in \{0,1\}} P(x_1 = v) P(x_2 = 1 \mid x_1 = v)$$

$$= 0.72$$

$$P(x_3 = 1) = \sum_{x_2 = v} P(x_2 = v) P(x_3 = 1 \mid x_2 = v)$$

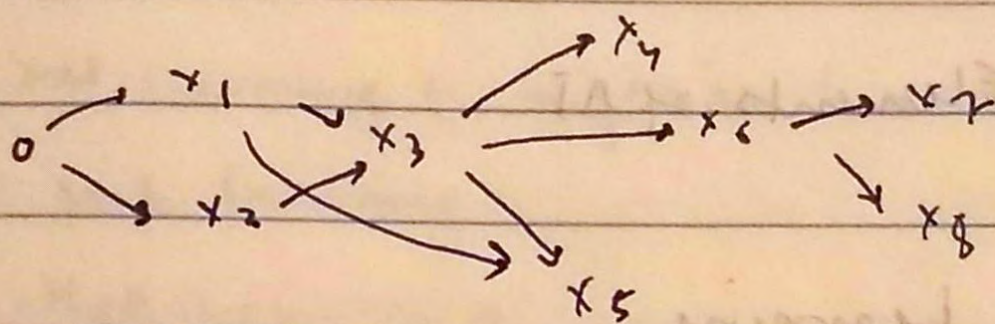
$$\begin{aligned}
 P(X_3=1) &= \sum_{x_1, x_2, x_4, x_5} P(x_1) P(x_2|x_1) P(x_3=1|x_2) P(x_4|x_3) P(x_5|x_4) \\
 &= \sum_{x_1, x_2, x_4} P(x_1) P(x_2|x_1) P(x_3=1|x_2) P(x_4|x_3) \sum_{x_5} P(x_5|x_4) \\
 &= \sum_{x_1, x_2} P(x_1) P(x_2|x_1) P(x_3=1|x_2) \sum_{x_4} P(x_4|x_3) \\
 &= \sum_{x_2} P(x_3=1|x_2) \sum_{x_1} P(x_1) P(x_2|x_1) \\
 &= P(x_3=1)
 \end{aligned}$$

$$(2) \quad P(X_1=1 | X_3=1) = \frac{P(X_1=1, X_3=1)}{P(X_3=1)}$$

$$(3) \quad P(X_3=1 | X_5=1) = \frac{P(X_3=1, X_5=1)}{P(X_5=1)}$$

$$P(X_3=1, X_5=1) = \sum_{x_1, x_2, x_4, x_5} P(x_1) P(x_2|x_1) P(x_3=1|x_2) P(x_4|x_3) P(x_5=1|x_4)$$

$$= \sum_{x_1, x_2} P(x_1) P(x_2|x_1) P(x_3=1|x_2) \sum_{x_4, x_5} P(x_4|x_3) P(x_5=1|x_4)$$

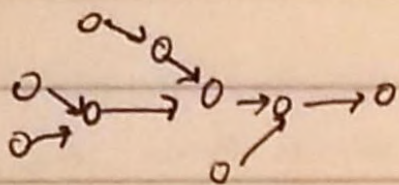


Sum x_7 \Rightarrow result depends on

$x_1, x_2, x_3, x_4, x_5, x_6$

Machine Learning

If Bayesian net is a polytree
then, var. elimination produces tables
which are \leq tables in original B.N.



if Bayesian net is not polytree,
tables could be $>$ than size of tables in B.N.

$$\sum_{x_i} f(x_1, x_2, x_3) g(x_1, x_5, x_6)$$

= function of (x_2, x_3, x_5, x_6)

for continuous R.V.,

$$P(x_1) = \int_{x_2} P(x_1 | x_2) P(x_2) dx$$

integrations should be simple

var. elim. is NP hard. for non polytrees

Belief Propagation (loopy B.P.)

assume graph is polytree
approximate inference.

Alternate idea:

instead of computing marginals exactly,
try to sample from the marginal
distribution.

*:

logistic regression

$P(w | \text{data})$

Previously, we computed a wrong

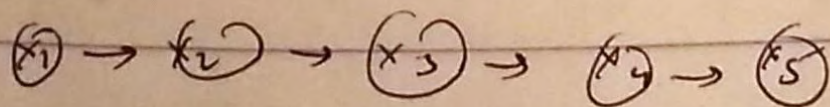
posterior $\hat{P}(w | \text{data})$

instead, try to sample w_1, w_2, \dots, w_k

$w_i \sim P(w | \text{data})$

- To predict instead of $\int P(w | \text{data}) P(y_{\text{next}} | w) dw$ pred dist.

compute $\frac{1}{k} \sum_i P(y_{\text{next}} | w_i)$



$P(x_2)$

$P(x_3 | x_1=1)$

$P(x_3 | x_5=1)$

sample x_1 , then $x_2 | x_1$, then $x_3 | x_2 \dots$

fix $x_1=1 \dots$

@ rejection sampling

Sampling \equiv monte-carlo.

① Rejection Sampling

take sample, discard if $x_5 \neq 1$
correct but slow (wastes sampling time)

② Likelihood weighting

instead of sampling and rejecting
stop after 1100 and force $x_5 = 1$
and use with weight of 0.3. $= \mathbb{E}P(x_5=1|x_4)$

11001 .3

10101 .3

00101 .3

0.9 total sample

does not waste samples. still slow.

For these algorithms, we must be able to
sample from $P(x_i | P_a(x_i))$

Easy for discrete. Various methods & tricks
for continuous variables.

Monte Carlo Markov Chain

Sample entire string at once.

$(101 \rightarrow) 10111 \rightarrow) 00101 \rightarrow) \dots$

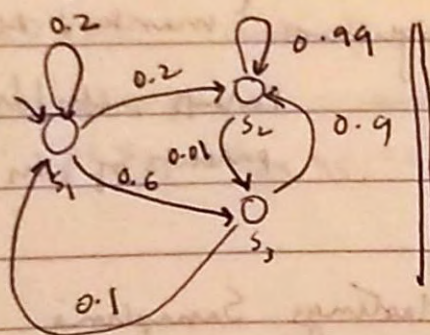
define a random process

in limit, the distribution is same as exactly what we want.

Markov Chain

Nodes are states.

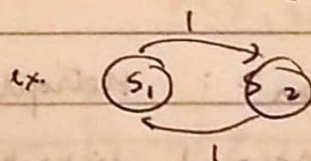
Edges are probabilistic transitions.



stationary distribution over states

$$\forall s \quad P(s) = P(\text{arrive in } s \text{ in next step})$$

$$\forall i \quad P(s_i) = \sum_j P(s_j) P(s_i | s_j)$$



does not have s.p.

We will build a Markov chain s.t. states are value configurations of Bayes net and its stationary distribution is $P(\text{unobserved vars.} | \text{observed vars.})$

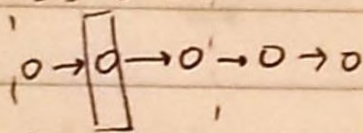
Transitions in the MC go from x_1, \dots, x_n to
(11001)

another valuation (01001)

Option 1: Gibbs Sampling

Pick $i \in \{1, \dots, N\}$ at random (uniformly)

Pick x_i from dist $P(x_i | x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$



need only neighbours "Markov blanket"

parents, children,
parents of children

Option 2: Metropolis's Hastings Sampling

* Proposal distribution $q(y|x)$

repeat

① Pick y from $q(y|x)$

② compute accept. probability $A = \min(1, \frac{P(y)q(x|y)}{P(x)q(y|x)})$

③ Accept y i.e. $x \leftarrow y$ with prob A

or stationary dist

Detailed Balance

Markov Chain with Transition prob T

has Detailed Balance relative to P

$$P(a) T(b|a) = P(b) T(a|b)$$

Fact 1 DB $\rightarrow P$ is stationary for MC

Fact 2 MH has DB for P .

Fact 3 Gibbs is special case of MH where $A=1$.

00100 current state

resample x_2 using Gibbs.

$$P(x_2 | x_{1,3,4,5} = 0100) = \frac{P(x_2 = v \text{ and } x_{1,3,4,5} = 0100)}{P(x_{1,3,4,5} = 0100)}$$

~~$$P(x_1=0) P(x_2=v | x_1=0) P(x_3=1) P(x_4=v)$$~~

numerator: $P(x_1=0) P(x_2=v | x_1=0) P(x_3=1 | x_2=v)$

$$P(x_4=0 | x_3=1) P(x_5=0 | x_4=0)$$

$$v = \{0, 1\}$$

$$P(x_2=0 \text{ and } x_{1,3,4,5} = 0100)$$

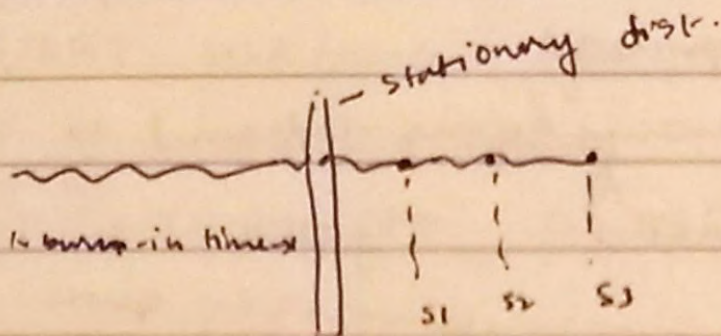
$$P(x_2=1 \text{ and } x_{1,3,4,5} = 0100)$$

$$P(x_2=0 | x_1=0) P(x_3=1 | x_2=0)$$

$$P(x_2=1 | x_1=0) P(x_3=1 | x_2=1)$$

Machine Learning

Markov Chain Monte Carlo.



s_1, s_2, s_3 are not independent
but we assume so.

Gibbs sampling

MH sampling

$$p(v_i | v_1, \dots, v_{i-1}, v_{i+1}, \dots, v_n) \quad \text{--- can't compute for undirected models}$$
$$\min \left(1, \frac{p(v') q(v|v')}{p(v) q(v'|v)} \right)$$

use MH

Fact 1: DB implies PC is stationary.

Fact 2: MH satisfies DB for PC.

Fact 3: Gibbs is MH & DB.

proof in slides

Latent Dirichlet Allocation

completely unsupervised text learning

Dirichlet

constraint $\sum \alpha_i = 1$

$$\alpha = (\alpha_1, \dots, \alpha_k)^T \quad \alpha_i = \sum d_i \quad \alpha_i - 1$$

$$p(w|d) = \text{Dir}(M|d) = \frac{\prod (d_{oi}) M_i}{\prod \Gamma(\alpha_i)}$$

Dirichlet list

$$p(\text{vector}) = \text{Dir}(M|k+m)$$

"topic" : dist over words.

each doc. is about multiple topics.

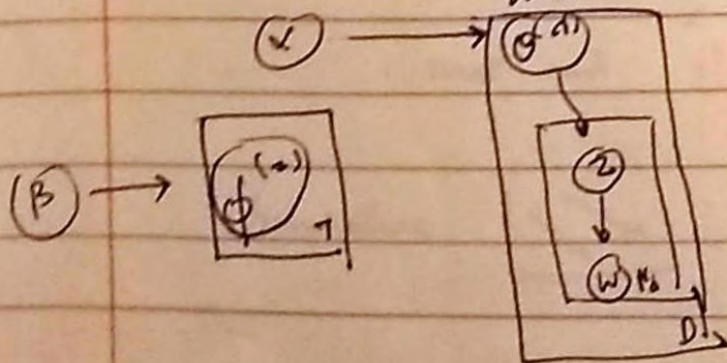
each doc. draws topic from dir.

each topic draws word from dir.

For each doc, for each iteration

decide topic

draw word from topic



LDA - matrix of doc/w. s.

Topic model - matrix of probs.

N : total # of words

N_d : # of words in doc d .

$N_{k,d}$: # of ~~times~~ times a topic occurred.

$N_{k,d}$: # of times topic k appeared in d .

$N_{i,k}$: # of times word i occurred in topic k .

prior: $\prod_{d \in D} \text{Dir}(\theta_d | \alpha)$

$\prod_{k=1}^K \text{Dir}(\phi_k | \beta)$

posterior: -

but we do not know topics (z_i)

Gibbs sampling

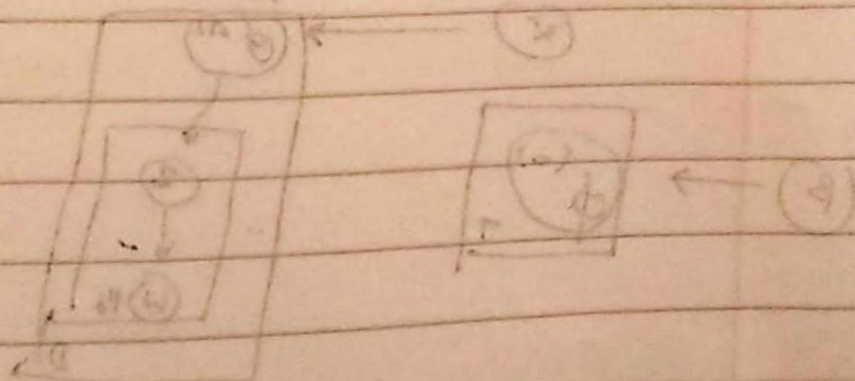
remove word.

Learning α and β .

Evidence Maximization

direct \rightarrow exponential

use sum of log of samples



Machine Learning

Y observed variables

Z hidden variables

want to max. $P(Y|\theta) = \int P(Y, Z|\theta) dz$

$$Q(\theta^{new}, \theta^{old}) = E_{P(Z|Y, \theta^{old})} [\ln P(Y, Z|\theta^{new})]$$

EM Algorithm

init θ^{old}

repeat

+ E step: calculate $Q(\theta^{new}, \theta^{old})$

* M step: pick θ^{new} to max $Q(\theta^{new}, \theta^{old})$

what does this mean?

Mixture of coin models

for each i

pick $z_i \in \{1, \dots, k\}$ from discrete $\{P_1, \dots, P_k\}$

ex. $k=3$ $z = \{1, 2, 3\}$

$p = \{0.7, 0.1, 0.2\}$

$z_i = 1 \iff 100$

$2 \iff 010$

$3 \iff 001$

flip coin z_i T times to get y_i

ex. $T=6$ $Y_i = UHTHTH$

coin i has $P(H) = \mu_i$

likelihood = $\mu_i^{\#H} (1-\mu_i)^{\#T}$

Complete data likelihood

$$L = \prod_i \prod_j \left[p_j \prod_t \mu_j^{y_{it}} (1-\mu_j)^{1-y_{it}} \right]^{z_{ij}}$$

coins \uparrow
 labels \uparrow

$$\theta = \{(p_j), \{\mu_j\}\}$$

$$\ln(P(Y, Z | \theta^{new})) = \sum_i \sum_j z_{ij} \left[\ln p_j + \sum_t y_{it} \ln \mu_j + (1-y_{it}) \ln (1-\mu_j) \right]$$

Next figure out \mathcal{Q} .

$$\mathcal{Q}(\theta^{new}, \theta^{old}) = \sum_i \sum_j E_{P(Z|Y, \theta^{old})} [z_{ij}] \dots$$

fn. of new
params

let $\alpha_{ij} = p_j P(Y_i | \theta_j) = p_j \prod_t \mu_j^{y_{it}} (1-\mu_j)^{1-y_{it}}$

fn. of old
params

let $\gamma_{ij} = E_{P(Z|Y, \theta^{old})} [z_{ij}] = \gamma_{ij}$

$$\mathcal{Q}(\theta^{new}, \theta^{old}) = \sum_i \sum_j \gamma_{ij} \ln \alpha_{ij}$$

$$Y_{ij} = E_{P(z|Y, \theta^{old})} [z_{ij}] = P(z_{ij}=1 | Y, \theta^{old})$$

$$= \frac{P(z_{ij}=1, y_i | Y^i, \theta^{old})}{P(y_i | Y^i, \theta^{old})}$$

$$= \frac{P(z_{ij}=1, y_i | \theta^{old})}{P(y_i | \theta^{old})} \quad \text{because } y_i \perp Y^{-i} | \theta^{old}$$

$$\text{Numerator} = P(z_i=j) P(y_i | z_i=j)$$

$$= p_j P(y_i | \theta_j)$$

$$= \underset{\text{binomial}}{\alpha_{ij}} \underset{\text{assumed}}{P(y_i | \theta_j)}$$

$$\Phi \left[Y_{ij} = \frac{\alpha_{ij}}{\sum_k \alpha_{ik}} \right] \quad \Gamma$$

$$Q = \sum_i \sum_j Y_{ij} \left[\ln p_j + \sum_t y_{it} \ln m_j + (1 - y_{it}) \ln (1 - m_j) \right]$$

p_j, m_j are new.

$$\text{Must satisfy } \sum_j p_j = 1 \quad \left\| \begin{array}{l} \text{use Lagrange} \\ \text{multipliers} \end{array} \right.$$

$$Q = \text{Const}(p_j) + \sum_i \sum_j Y_{ij} \ln p_j$$

$$L = \sum_i \sum_j Y_{ij} \ln p_j + \lambda (\sum_j p_j - 1)$$

No. of examples in class j

$$\frac{\partial}{\partial \lambda} \sum p_j - 1 = 0; \quad \frac{\partial}{\partial p_j} \sum_i \sum_t r_{ij} \cdot \frac{1}{p_j} + \lambda = 0$$

$$\Rightarrow p_j = -\frac{1}{\lambda} \sum_i r_{ij}$$

$$= -\frac{1}{\lambda} N_j$$

$$N_j = \sum_i r_{ij}$$

$$\sum_j p_j = 1 \Rightarrow \sum_j -\frac{1}{\lambda} N_j = 1 \Rightarrow N = \lambda$$

$$p_j = \frac{N_j}{N} \quad \text{II}$$

$$\frac{\partial}{\partial \mu_j} = \sum_i \sum_t \left[\frac{y_{it}}{\mu_j} - \frac{1-y_{it}}{1-\mu_j} \right] r_{ij} = 0$$

$$\sum_i \sum_t (y_{it}(1-\mu_j) - \mu_j(1-y_{it})) r_{ij} = 0$$

$$\sum_i \sum_t r_{ij} (y_{it} - y_{it} \mu_j - \mu_j + \mu_j y_{it}) = 0$$

$$\sum_i \sum_t r_{ij} \mu_j = \sum_i \sum_t r_{ij} y_{it}$$

$$\Rightarrow \mu_j \sum_i \sum_t r_{ij} = \sum_i r_{ij} \left(\sum_t y_{it} \right)$$

$$\Rightarrow \mu_j T N_j = \sum_i r_{ij} \left(\sum_t y_{it} \right)$$

$$\Rightarrow \mu_j = \frac{\sum_i r_{ij} \left(\sum_t y_{it} \right)}{T N_j} \quad \text{III}$$

EM for Mixture of Coins

init $\{p_j, u_j\}$

Repeat

Calculate T_{ij} using I

Calculate $\{p_j, u_j\}$ using II, III

EM does marginal likelihood estimation when one or more parameters are not observed.

Why does it work?

Fact if $Q(\theta^{new}, \theta^{old}) > Q(\theta^{old}, \theta^{old})$

then $P(Y|\theta^{new}) > P(Y|\theta^{old})$

Proof $0 < Q(\theta^{new}, \theta^{old}) - Q(\theta^{old}, \theta^{old})$

$$= E_{P(z|Y, \theta^{old})} [\ln P(Y, z|\theta^{new}) - \ln P(Y, z|\theta^{old})]$$

$$= E_{P(z|Y, \theta^{old})} \left[\ln \frac{P(Y|\theta^{new}) P(z|Y, \theta^{new})}{P(Y|\theta^{old}) P(z|Y, \theta^{old})} \right]$$

$$= E_{P(z)} \left[\ln \frac{P(Y|\theta^{new})}{P(Y|\theta^{old})} + E_{P(z)} \left[\ln \frac{P(z|Y, \theta^{new})}{P(z|Y, \theta^{old})} \right] \right]$$

no 2 term

less than 0
claim

$$0 \leq \ln \frac{P(Y|\theta^{new})}{P(Y|\theta^{old})} \Rightarrow \frac{P(Y|\theta^{new})}{P(Y|\theta^{old})} \geq 1$$

Jensen's inequality: concave f

ex. variance

$$E[X]^2 - E[X^2] \geq 0$$



$$f(E[X]) \geq E[f(X)]$$

$P_1(v)$ $P_2(v)$

k_L divergence
not symmetric

$$d_{k_L}(P_1 \parallel P_2) = \int P_1(v) \ln \frac{P_1(v)}{P_2(v)} dv$$

$$= E_{P_1(v)} \left[\ln \frac{P_1(v)}{P_2(v)} \right] \geq 0$$

$$-d_{k_L}(P_1 \parallel P_2) = E_{P_1(v)} \ln \frac{P_2(v)}{P_1(v)}$$

$$\leq \ln E_{P_1(v)} \frac{P_2(v)}{P_1(v)}$$

$$= \ln \int P_1(v) \frac{P_2(v)}{P_1(v)} dv = \ln 1 = 0$$

Machine Learning

Kernel Function : Fast way to compute inner product in some feature space

$$k(x_i, x_j) = \phi(x_i)^T \phi(x_j)$$

For ex. euclidean $\frac{(x_i^T x_j)}{\frac{1}{2} \|x_i - x_j\|^2}$, quadratic $(x_i^T x_j + 1)^2$, polynomial $(x_i^T x_j + 1)^d$

int. polynomial \rightarrow RBF (e)

Alg \leftarrow kernel \leftarrow para

Kernel Method : Learning algorithm that works with kernels

For ex. Perceptron, KNN, Regularized Linear Regression.

Merker's Theorem : $k(\cdot, \cdot)$ is a kernel \Leftrightarrow

- ① k is symmetric ② kernel matrix is Positive Semidefinite \forall finite points

All eigenvalues $\geq 0 \Leftrightarrow \forall c \ c^T K c \geq 0$

Proof one dim. if $k(a, b) = \phi(a)^T \phi(b)$ then $c^T K c \geq 0$

$$\begin{aligned} c^T K c &= \sum_i \sum_j c_i c_j k(x_i, x_j) \\ &= \sum_i \sum_j c_i c_j \phi(x_i)^T \phi(x_j) \\ &= \left[\sum_i c_i \phi(x_i) \right]^T \left[\sum_j c_j \phi(x_j) \right] \\ &= \left\| \sum_i c_i \phi(x_i) \right\|^2 \geq 0 \end{aligned}$$

Fact if k_1 is a kernel, k_2 is a kernel

then ① $k_3 = k_1 + k_2$ is a kernel

② $k_4 = k_1 \times k_2$ is a kernel.

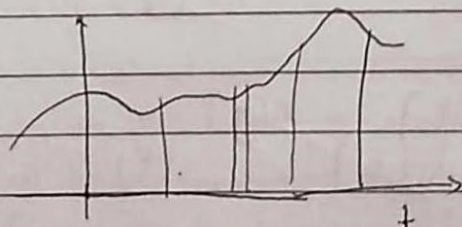
Proof of ① $c^T k_3 c = c^T (k_1 + k_2) c = c^T k_1 c + c^T k_2 c \geq 0$

After $k_3(x_i, x_j) = k_1(x_i, x_j) + k_2(x_i, x_j)$

$$\begin{aligned} &= \phi_1(x_i)^T \phi_1(x_j) + \phi_2(x_i)^T \phi_2(x_j) \\ &= \hat{\phi}(a)^T \hat{\phi}(b) \quad \left| \begin{array}{l} \hat{\phi}(a) = \text{concat} \ \phi_1(a) \\ \phi_2(a) \end{array} \right. \end{aligned}$$

Gaussian Process

A distribution over functions such that for any finite set of inputs x_1, \dots, x_N the vector of function values



$f = (f(x_1), \dots, f(x_N))^T \sim N((m(x_1), \dots, m(x_N))^T, C)$

is distributed normally

it is specified by mean fn. and covariance

$$C_{ij} = \begin{pmatrix} & j \\ i & \end{pmatrix} = k(x_i, x_j)$$

next point is first in
covariance matrix

Given $x_1 \dots x_N$
 $f_1 \dots f_N$
 x_{N+1}

$$C_{N+1} = \begin{pmatrix} C & v^T \\ v & C_N \end{pmatrix} \quad \begin{matrix} 0 = C(x_{N+1}, x_{N+1}) \\ v = C(x_1, x_{N+1}) \\ \vdots \\ C(x_N, x_{N+1}) \end{matrix}$$

Predict f_{N+1}

$$\bar{f}_{N+1} = (f_1 \dots f_N) \quad C_N = C \text{ applied to } x_1 \dots x_N$$

Assume $m(x) = 0 \quad \forall x$

$$\bar{f}_{N+1} = (f_{N+1} \quad (\bar{f}_N)^T)^T \sim N(0, C_{N+1})$$

observe \bar{f}_N

$p(f_{N+1})$

$$\Sigma = C - v^T C_N^{-1} v$$

$$M = 0 + v^T C_N^{-1} (\bar{f}_N)$$

Application
of MVN
conditional
template

$$\begin{pmatrix} x_a \\ x_b \end{pmatrix} \sim N \left(\begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix}, \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix} \right)$$

$$\Sigma_{a|b} = \Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba}$$

$$M_{a|b} = \mu_a + \Sigma_{ab} \Sigma_{bb}^{-1} (x_b - \mu_b)$$

Bayes LR.

$$p(w) = N(0, \frac{1}{\alpha} I) ; y = \Phi w ; t \sim N(y, \frac{1}{\beta} I)$$

$$t_i \sim N(w^T \phi(x_i), \frac{1}{\beta})$$

(Claim: t is sampled from gaussian process.

$$E[y] = \Phi E[w] = 0$$

$$E[y y^T] = E[\Phi w w^T \Phi^T]$$

$$= \Phi E[w w^T] \Phi^T$$

$$= \Phi \left[\frac{1}{\alpha} I \right] \Phi^T = \frac{1}{\alpha} \Phi \Phi^T$$

$$\Phi = \begin{pmatrix} \phi(x_1) \\ \vdots \\ \phi(x_N) \end{pmatrix}_{N \times D}$$

$$\Phi^T \Phi = \begin{pmatrix} \phi(x_1)^T \\ \vdots \\ \phi(x_N)^T \end{pmatrix} \begin{pmatrix} \phi(x_1) \\ \vdots \\ \phi(x_N) \end{pmatrix}$$

$$= K \quad N \times N$$

$$E[t] = E[y] = 0$$

$$E[t t^T] = \text{cov}(y) + \text{cov}(t|y)$$

$$= \frac{1}{\alpha} K + \frac{1}{\beta} I$$

$$\Phi \Phi^T = \begin{pmatrix} \phi(x_1) \phi(x_1)^T \\ \vdots \\ \phi(x_N) \phi(x_N)^T \end{pmatrix}$$

inner product

For any $x_1 \dots x_N$

$$t = (t(x_1) \dots t(x_N)) \sim N(0, C)$$

$$C = \frac{1}{\alpha} K + \frac{1}{\beta} I$$

$$C = \frac{1}{\alpha} K + \frac{1}{\beta} I$$

$$K_{ij} = e^{-\frac{1}{2} \|x_i - x_j\|^2 / s^2}$$

$$\text{Evidence} = P(t | \alpha, \beta) = N(0, C_N)$$

$$\log \tilde{E} = -\frac{N}{2} \log 2\pi - \frac{1}{2} \log |C_N| - \frac{1}{2} t^T C_N^{-1} t$$

= Marginal Likelihood.

$$\frac{\partial}{\partial \alpha} \quad \frac{\partial}{\partial \beta} \quad \frac{\partial}{\partial s} \quad \Rightarrow \text{Gradient Descent}$$

Logistic - Bernoulli?
if $m \neq 0$ then?