

• are trademarks, real numbers

Supervised learning

```

graph TD
    TS[Training set] --> LA[Learning algorithm]
    LA --> TS2[Test set]
    TS2 -- "(\"output\")" --> M[M]
    TS2 -- "estimated output by e.g. f(x) ..." --> EOE[...]
  
```

$$L(x) = \sum_{i=1}^n \theta_i x_i + \theta_0$$

for concavity, define $L(x)$

$$L(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

x_1, x_2 = number of features
 $\theta_0, \theta_1, \theta_2$ = coefficients

$$L(x) = \frac{1}{n} \sum_{i=1}^n L_i$$

this is loss function

$$\begin{aligned} x_1 &= \text{no. of books} \\ x_2 &= \text{no. of pens} \end{aligned}$$

$x_1 + x_2 = \text{total no. of items}$

The diagram illustrates a supervised learning process. At the bottom, a box labeled "Training set" has an arrow pointing up to a box labeled "(Learning algorithm)". From this middle box, an arrow points up to a box labeled "Test set". From the "Test set" box, an arrow points left to a box labeled "('Input features')". From this input box, an arrow points left to a box labeled "Estimated output". The word "New data" is written vertically to the right of the "Test set" box.

$$\begin{array}{|c|} \hline \text{In } X_{\perp -}(X_{\perp}X) = 0. \\ \hline \end{array}$$

Indoor & outdoor

~~26.N~~

12 / sec ; digitize image, & feed to NN!
After 2 mins, it learns to identify the neutrinos
of the human driver! repeat for other road types!
Sueanne at Carnegie Mellon, 1996

If we car driving itself - never need drivers!

$$x_1 - (x_2 x) = \theta.$$

residual problem: solving as continuous value

m = number of training examples
 x = "input" variable / "features"
 y = "output" variable / "target" variable
 n = \sum training examples : $(x^{(1)}, y^{(1)})$
 \dots
 $(x^{(n)}, y^{(n)})$

Notes

$$\text{for } f_{\theta}(\theta) : \quad \hat{\theta} = \arg \min_{\theta} \sum_{i=1}^n \frac{1}{n} (y_i - (\theta x_i))^2$$

to find parameters

Input variable example: (x_1, x_2, \dots, x_n) = I feature example
Output variable example: y = "output variable / "target", variable

Batch: lots of whole training set!

- steeper descent!
- the derivative gives it to you ... that is the direction of

How is gradient descent (local) algorithm?

(θ) $J(\theta)$
 ("") is doesn't change
How to tell for convergence?



As you approach minimum, decrease the gradient $\rightarrow 0$,
 so smaller steps



for OLS, $J(\theta)$ is quadratic, & a nice bowl

Repeating it! (Convergence)

$$\theta_i := \theta_i - \alpha \sum_j x_{ij} (h_\theta(x_j) - y_j) \cdot x_{ij}$$

Batch gradient descent

for m training examples

α - alpha - the learning rate! size of the update:

$$\theta_i = \theta_i - \alpha (h_\theta(x_i) - y_i) \cdot x_i$$

out the matrix

$$= (h_\theta(x) - y) \cdot \frac{\partial}{\partial \theta} (\theta_0 x_0 + \dots + \theta_n x_n - y) \text{ rule: chain rule}$$

apply:

$$\frac{\partial}{\partial \theta_i} J(\theta) = \frac{\partial}{\partial \theta_i} (h_\theta(x) - y)^2$$

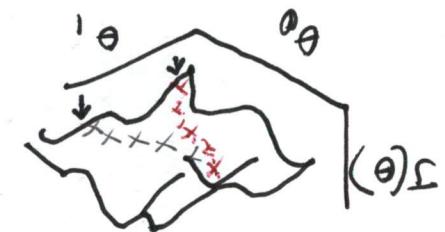
for 1 training example

$a = b$: means
an algorithm finds
the equal, in
a variable to
"self"

$$\theta_i := \theta_i - \alpha \frac{\partial}{\partial \theta_i} J(\theta)$$

where you end up can depend on where you

step to go down fastest
look around, small
step at a time



keep changing θ , to minimize $J(\theta)$

start at same θ ($\theta = 0$: vector of 0s)

[47:40]



minimum! keep near boundary a bit
much faster, but won't converge to global
quickly

last at 15+ training epochs; difficult to note

} { for all?

$\theta_i := \theta_i - \alpha (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$

$f_{\theta_i} = 1 + w_i^T$

repeat until convergence {

for large training sets: stochastic gradient descent

[Ng 2]

(3)

linear vector

$$\Delta \theta = \theta$$

gradient descent

linear

gradient descent

$$A = \sum_{i=1}^m A_i$$

$$B = \sum_{i=1}^m B_i$$

③ top

linear

$$\begin{bmatrix} \frac{\partial C}{\partial C} \\ \frac{\partial C}{\partial C} \\ \vdots \\ \frac{\partial C}{\partial C} \end{bmatrix} = (\nabla f)^T \Delta$$

update

$$A = \sum_{i=1}^m A_i$$

$$\nabla f = \sum_{i=1}^m \nabla f_i$$

1 - 2

gradient descent

numerical learning

$\times \times \times \times \times \times$

\$\\$

linear regression

function 2

No

1 + m

No update

first

numerical
learning

gradient

① top

$$\begin{bmatrix} \frac{\partial C}{\partial C} \\ \frac{\partial C}{\partial C} \\ \vdots \\ \frac{\partial C}{\partial C} \end{bmatrix}$$

②

① top

$$\begin{aligned}
 & \text{Left side: } h - (\theta^T x) y = h - \theta^T x - \theta^T x y \\
 & \text{Right side: } \frac{1}{n} \sum_{i=1}^n [h_i - \theta^T x_i] y_i = \frac{1}{n} \sum_{i=1}^n [h_i - \theta^T x_i + \theta^T x_i y_i] = \frac{1}{n} \sum_{i=1}^n h_i + \frac{1}{n} \sum_{i=1}^n \theta^T x_i y_i \\
 & \quad \text{Since } \frac{1}{n} \sum_{i=1}^n h_i = h \text{ and } \frac{1}{n} \sum_{i=1}^n \theta^T x_i y_i = \theta^T \left(\frac{1}{n} \sum_{i=1}^n x_i y_i \right) = \theta^T \bar{x} y \\
 & \quad \text{Therefore, } h - (\theta^T x) y = h - \theta^T x
 \end{aligned}$$

$$\begin{aligned}
 & \text{Left side: } \theta \cdot u = \text{feature vector} \\
 & \text{Right side: } w \cdot v = \text{feature vector} \\
 & \text{Left side: } \begin{bmatrix} (w^T x)^{\theta_1} \\ \vdots \\ (w^T x)^{\theta_m} \end{bmatrix} = \begin{bmatrix} \theta_1(w^T x) \\ \vdots \\ \theta_m(w^T x) \end{bmatrix} = \theta X \\
 & \text{Right side: } \theta \left[\begin{array}{c} \hline \quad \downarrow (w^T x) \\ \hline \quad \downarrow (w^T x) \\ \hline \quad \downarrow (w^T x) \end{array} \right] = \theta X
 \end{aligned}$$

$$A + ABA^T C = CAB + C^T A B^T$$

$$\text{Tr } ABL = \text{Tr } \cancel{C}AB = \text{Tr } (BCA)$$

Ng 2-3

Normal eqn $\left[\begin{matrix} h_1^T X = \theta X^T X \end{matrix} \right]$

so \Rightarrow closed form

$$h_1 X_{T-1}(X^T X) = \theta$$

θ (left gradient descent)

$$\theta_j := \theta_j + \frac{\partial}{\partial \theta_j} L(\theta) = \theta_j - \frac{1}{n} \sum_{i=1}^n (y_i - h_\theta(x_i)) x_i$$

compute partial deriv for each w.r.t θ_j :
 algebra

$$\theta := \theta + \Delta \theta (\theta) : \text{note the quadratic}$$

apply gradient descent algorithm \leftarrow gradient descent

$$(x_i^\theta y_i - 1) \log(y_i) + (1 - y_i) \log(1 - y_i) = \theta^\top x_i$$

$$\text{find } \theta \text{ to fit max } L(\theta)$$

$$L(\theta) = P(y|x) = P(y|x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(y_i|x_i)$$

(if $y=0$, one left disappears!):
 compute weight of

$$P(y|x^\theta y - 1) = h_\theta(x^\theta y) = P(y|x)$$

$$P(y|x^\theta y - 1) = h_\theta(x^\theta y) = h_\theta(x) \quad \text{softmax function}$$

$$h_\theta(x) = g(\theta^\top x) = \frac{1}{1 + e^{-\theta^\top x}}$$

$$h(x) \in [0, 1] \quad y \in \{0, 1\}$$

(2)

add power, the line class

~~$$\theta_j := \theta_j - \frac{\partial}{\partial \theta_j} L(\theta) = \theta_j - \frac{1}{n} \sum_{i=1}^n (y_i - h_\theta(x_i)) x_i$$~~

not a good idea
 apply CR to column = $\theta \cdot k_y$ but

will the computation crash
 in the small space?

will the loss oscillate!

gradient problem: do tiny move to discuss
 classification: $y \in \{0, 1\}$

endows L with probability properties

θ does not matter: — you get same values of θ !

$$L(\theta) = \sum_{i=1}^n (y_i - \theta^\top x_i)^2$$

so max $L(\theta)$ is the same as minimum

$$\left[\left(\frac{x_i}{\theta} \right)^2 + \left(\frac{y_i}{\theta} \right)^2 \right] = \left[\frac{x_i^2}{\theta^2} + \frac{y_i^2}{\theta^2} \right] = \frac{1}{\theta^2} (x_i^2 + y_i^2) = \frac{1}{\theta^2} \exp(-\dots) = \exp \left(\frac{-x_i^2 - y_i^2}{\theta^2} \right)$$

$$L(\theta) = \log L(\theta)$$

choose θ to maximize $L(\theta) = P(y|x^\theta)$ or
 probably no possibility to see dot

maximum likelihood:

$$Ng 3$$

Ng3

- simpler than logistic regression
- difficult to extend to probabilistic softmax

Curry - $\{x \cdot ((\lambda - h(x)) + !\theta) = 0\}$

$h(x) = (x, \theta)^T$

$g(z) = \begin{cases} 0 \\ 1 \end{cases}$

Percptron algorithm



To learn: record if! carry-up the derive! do it

Some learning rule / algorithm

$h_\theta(x)$ is not the answer, but logistic function!

Is it the same algorithm? No. In logistic regression

(3)

$$(n) \frac{d}{dx} = \frac{d}{dx}$$

$$\int_0^x (n-1) \frac{d}{dt} dt + \log(1-y) = \exp(\log(1-y))$$

$$(n-1)(y-1) + \log(y) = \exp(y)$$

$$(n-1)(y-1) \frac{dy}{dx} = \exp(y)$$

$$n-1 = (\frac{dy}{dx})^{-1} = \frac{1}{\exp(y)} = \frac{1}{y^n}$$

for a, b, t , we can get different results.

$$y = f(t) \text{ (function of time)}$$

$$y = f(t) = \frac{1}{\exp(-t)} = \frac{1}{e^{-t}} = e^t$$

$$\frac{dy}{dt} = \frac{d}{dt} \left(\frac{1}{e^t} \right) = -e^{-t} = -\frac{1}{e^t}$$

$$\frac{d}{dt} \left(\frac{1}{e^t} \right) = -e^{-t} = -\frac{1}{e^t}$$

on every measure: if $t=0$, you have to invert it

$$\text{where } H = \text{Hessian matrix: } -\frac{1}{2} \text{ 2nd derivative}$$

$$\theta^{t+1} = \theta^t - H^{-1} \Delta \theta$$

N.M for θ as a vector

⑥

Newton M step for convexity (quadratic convergence)

All of these algorithms converge and the speeds are different!

All of these algorithms converge and the speeds are

$$\theta^{t+1} = \theta^t - \frac{\nabla J(\theta^t)}{J''(\theta^t)} \quad \text{of gradients for loss regularization}$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \ell(\theta) + \lambda R(\theta)$$

$$\theta^{t+1} = \theta^t - \frac{\nabla J(\theta^t)}{J''(\theta^t)}$$

$$\theta^{t+1} = \theta^t - \frac{\nabla J(\theta^t)}{J''(\theta^t)}$$

$$\theta^{t+1} = \nabla \left| \frac{\nabla J(\theta^t)}{J''(\theta^t)} \right|$$

$$\theta^{t+1} = \theta^t - \frac{\nabla J(\theta^t)}{J''(\theta^t)}$$

$$\theta^{t+1} = \theta^t - \frac{\nabla J(\theta^t)}{J''(\theta^t)}$$

Newton's method is faster

to fit together batch/stochastic gradient descent is slow.

$$\theta^{t+1} = \theta^t + \alpha (y^{(t)} - h(x^{(t)})) x^{(t)} \rightarrow \text{stoch grad descent}$$

$$(y^{(t)} - h(x^{(t)})) x^{(t)} + \frac{1}{m} \sum_{i=1}^m (y^{(i)} - h(x^{(i)})) x^{(i)} = J'(\theta)$$

$$\theta^{t+1} = \theta^t + \frac{1}{m} \sum_{i=1}^m (y^{(i)} - h(x^{(i)})) x^{(i)}$$

generalized linear models

2. explicit family

fitting. Logs \rightarrow in addition we need

[Nq4]

$$(n) \text{ Using } x = \frac{\theta!x^h}{\theta^h h!}$$

$$\text{run of } m \text{ groups : } \sum_{\perp} \theta = n \quad (3)$$

3

$$\begin{aligned}
 & 1 \rightarrow \emptyset \dots \emptyset \emptyset : \emptyset \dots \\
 & (1 \rightarrow \emptyset + \dots + 1 \emptyset) - 1 = \emptyset \\
 & P(y=i) = \emptyset \\
 & \text{Parameters: } \emptyset, \emptyset, \dots, \emptyset \\
 & y \in \{1, 2\} =
 \end{aligned}$$

Multinomial: ~~example~~ (!!!)

Gegevens Gaussian exponetieel: de som van i^k (combinaties) \rightarrow $\sum_{i=0}^{n-1} i^k = \frac{1}{k+1} [n : (k+1)]_E = \binom{n}{k} k!$

$$\log \text{residuals} = \frac{\ln \frac{x}{x_0} + c}{-6T} = \frac{\ln \frac{x^2+1}{1}}{1} = \frac{\ln x^2 + \ln 1 + 1}{1} =$$

$$(g!x) \tau = h \circ [g!x]^{h^{-1}} \exists = (x \# h)$$

g!x \in \text{defn of } g

$$(iv) \text{Growth} \propto \frac{\text{Exponential}}{\text{Exponential}}$$

$$\text{! univerzalni suvremeni } \alpha = [i:n] \quad x_1 \alpha = n (3)$$

How many children are there in your family? How do you divide a bill?

Most effective dists can be written in terms of exp. form

$$e.g. \text{multinat } N! \text{ multinomial! process ... costs}$$

$$\text{or Gaussian } (\geq 0)! \text{ exponential } (\geq 0)! \text{ beta! difficult features} \\ = \text{dist + extra features}$$

$$B = (M_1 \perp \boxed{M_2 = M_1}) (M_2 \text{ free})$$

$$(2\pi l^2 \frac{1}{2} M_2^2 + h^2) dx \left(z^{l^2 \frac{1}{2}} - \right) ds \frac{dx}{T} =$$

$\frac{N(\mu, \sigma^2)}{\text{for } \text{variance}}$: $\text{set } \sigma^2 = 1$; μ is a scaling factor

remove ref fpd

$$\text{Top} \leftarrow \left\{ (k+1)S_{0j} = (k-1)S_{0j} - \frac{1}{k-1} = (k)S_{0j} \right.$$

$$\frac{b-1}{b} \cdot b = n$$

(两边约分)

(5)

$$\left\{ \begin{array}{l} \text{...} \\ \text{...} \\ \text{...} \end{array} \right. \quad \left. \begin{array}{l} x_{1,0}^{(1)} \frac{\sum_{i=1}^{n+1}}{x_{1,0}^{(0)}} = 1 \\ x_{1,0}^{(1)} \\ \vdots \\ x_{1,n}^{(1)} \frac{n!}{n!} = 1 \end{array} \right. \quad \text{so } \frac{(x_1)^{n+1}}{n!} = 1$$

$$\left(\theta! x | (h) \right) \cdots \left(h | (x) \right) \text{ für } \begin{array}{l} \text{die Summe der} \\ \text{Ziffern ist gleich der} \\ \text{Summe der Ziffern im Produkt!} \end{array}$$

$$\left[\begin{array}{c} \left(x_{1,0}^{(1)} \frac{\sum_{i=1}^{n+1}}{x_{1,0}^{(0)}} \right) / x_{1,0}^{(1)} \\ \vdots \\ \left(x_{1,n}^{(1)} \frac{n!}{n!} \right) / x_{1,n}^{(1)} \end{array} \right] \xrightarrow{\text{so } \frac{(x_1)^{n+1}}{n!} = 1}$$

$$\left[\begin{array}{c} 1 \\ \vdots \\ 1 \end{array} \right] = \left[\begin{array}{c} \theta! x | (h) \\ \vdots \\ 1 = h \end{array} \right] \quad \text{so LAGRANGE: } \frac{\theta! x | (h)}{y} = (x)^{\theta}$$

(6)

$$\left[\begin{array}{c} x_{1,0}^{(1)} = ? \\ \vdots \\ x_{1,n}^{(1)} = ? \end{array} \right] \quad \frac{x_{1,0}^{(1)} \sum_{i=1}^{n+1}}{x_{1,0}^{(0)}} + 1 =$$

$$(1 - \dots - 1) \cdot \frac{1}{1 - \frac{1}{x_{1,0}^{(1)}}} = 1$$

map
der
zu \mathbb{R}^k

$$\begin{cases} \text{...} \\ \text{...} \\ \text{...} \end{cases} \quad \begin{cases} T_F(h) : (1 - \dots - 1) \cdot \frac{1}{1 - \frac{1}{x_{1,0}^{(1)}}} = 1 \\ \vdots \\ ((h) - (h, 1, \dots, 1)) \cdot \frac{1}{1 - \frac{1}{x_{1,0}^{(1)}}} = 1 \end{cases} \quad \text{so } \frac{1}{1 - \frac{1}{x_{1,0}^{(1)}}} = 1$$

$$\left[\begin{array}{c} 1 \\ \vdots \\ 1 \end{array} \right] = \left[\begin{array}{c} 1 \\ \vdots \\ 1 \end{array} \right] \quad \text{so } \frac{1}{1 - \frac{1}{x_{1,0}^{(1)}}} = 1$$

$$\left[\begin{array}{c} 1 \\ \vdots \\ 1 \end{array} \right] = \left[\begin{array}{c} 1 \\ \vdots \\ 1 \end{array} \right] \quad \text{so } \frac{1}{1 - \frac{1}{x_{1,0}^{(1)}}} = 1$$

$$T(I) = (I) \quad T(0) = (0) \quad (1 - \dots - 1) \quad \text{K-1 dimensional vector}$$

$$T((1, 2, \dots, k)) = (1, 2, \dots, k) \quad \text{so } \frac{1}{1 - \frac{1}{x_{1,0}^{(1)}}} = 1$$

Multinomial cat

Nag 4 (cont.)

$$\max_{\phi} \mathcal{L}(\phi, w_0, u_i, z) = \frac{\sum_{i=1}^n \left(y_i - \sum_{j=1}^m w_j x_{ij} \right)^2}{\sum_{i=1}^n \left(y_i - \sum_{j=1}^m w_j x_{ij} \right)}$$

$$\text{MLE} \rightarrow \text{Conditional likelihood}$$

$$\left(\sum_{i=1}^n h_i(x) \right)^{\phi} = \prod_{i=1}^n h_i^{\phi}(x) = \prod_{i=1}^n \left(\sum_{j=1}^m a_{ij} x_j \right)^{\phi} = \prod_{i=1}^n \left(\sum_{j=1}^m a_{ij} \sum_{k=1}^n b_{kj} x_k \right)^{\phi} = \prod_{i=1}^n \left(\sum_{k=1}^n \sum_{j=1}^m a_{ij} b_{kj} x_k \right)^{\phi} = \prod_{k=1}^n \left(\sum_{i=1}^n \sum_{j=1}^m a_{ij} b_{kj} \right)^{\phi} x_k^{\phi} = \prod_{k=1}^n P(x_k)^{\phi} = P(x)^{\phi}$$

1

Useful tool at learning on the net.

$\begin{bmatrix} 1 & 5 & -5 \\ 1 & 1 & -5 \end{bmatrix}$

$\begin{bmatrix} 0 & 3 & 0 \\ 1 & 1 & 0 \end{bmatrix}$

$\begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \end{bmatrix} = 3$

$P(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x-\mu)^2\right)$ - multivariate normal

$\bar{x} \sim N(\mu, \Sigma)$... ~~covariance matrix~~ $Ex - \mu^T$

Assume $P(x|y)$ is known

column-wise discrimination always algorithm

Assume $x \in \mathbb{R}^n$, unknown

Example

where $P(x) = P(y=0|x)P(x) + P(y=1|x)$

$P(x)$

$P(y=1|x) = \frac{P(x|y=1) \cdot P(y)}{P(x)}$

builds a probability model of features conditioned as class labels y .

$P(y|x) = \text{Learn}$

$P(y|x) = \text{Learn}$

$\boxed{\text{This is inductive}}$

$\boxed{\text{Learn}}$

Naive Bayes

Classification

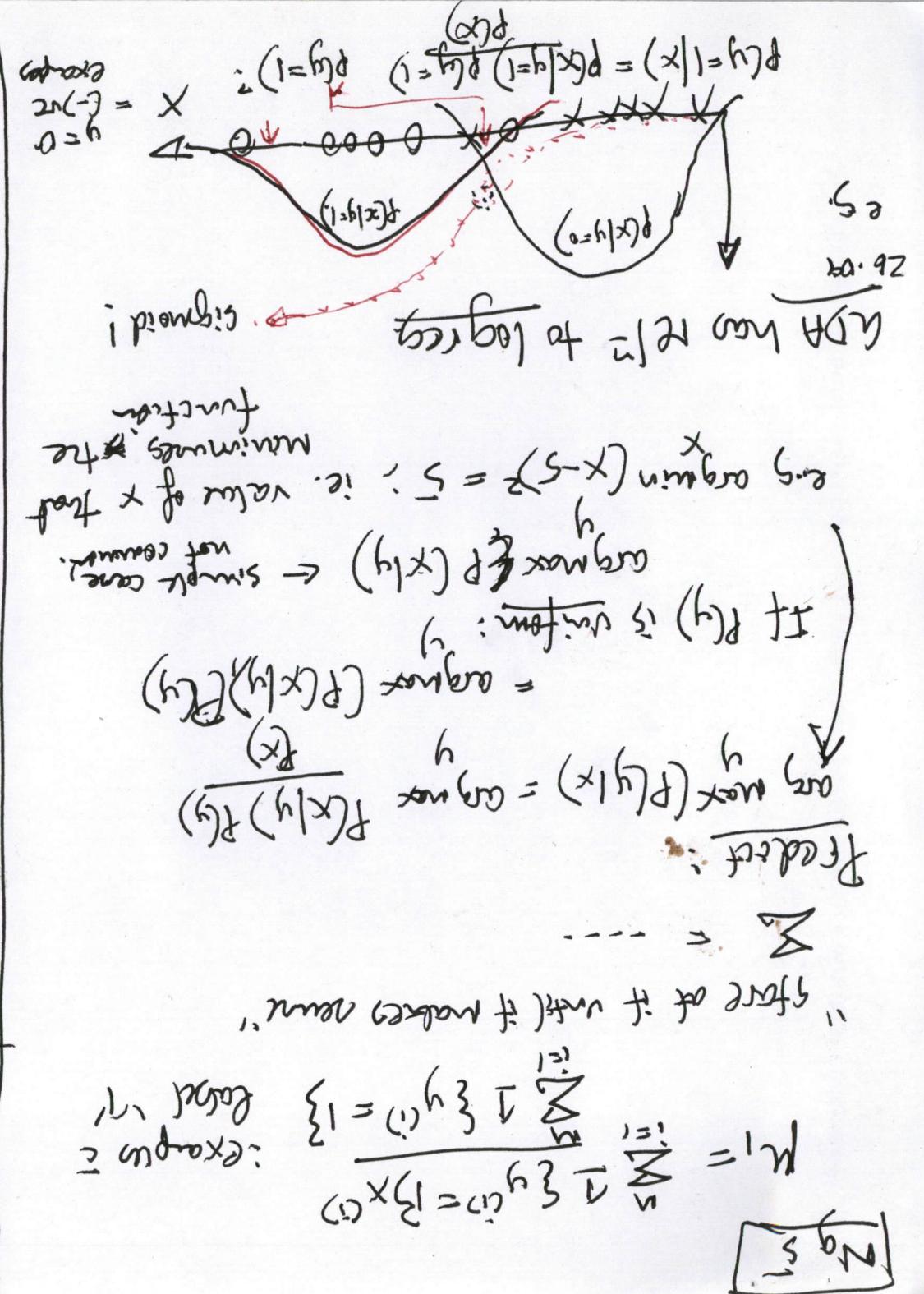
SLK View

CDH

N 9 5

- (41.30)
- High dimension shows artifacts
 - Probability distribution of ϕ vs HNE
 - training data
 - log loss needs less computation, but needs more data
 - GDA needs less training data
 - Real advantage of gda classifier: require less data!
 - will do better.
 - But if you don't know shape of $P(y|x)$, then log loss
 - it makes better use of the data
 - when $x|y$ is Gaussian, then GDA will do better
 - (real-life) DNN GDA & log loss
 - shows robustness of log loss to model assumptions
 - e.g. $y = \log(x)$ is linear + iid
 - $x|y \sim \text{Poisson}(x_0)$. Gaussian assumption
 - $\log(x) \sim \text{Poisson}(\lambda_0)$. Exponential assumption
 - $\log(x) \sim \text{Exponential}(\lambda_0)$. Poisson assumption
 - $\log(x) \sim \text{Normal}(\mu_0, \sigma^2)$. Gaussian assumption
 - logistic posterior for $\phi(y|x)$
 - $x|y \sim \text{Gaussian}$
-
- Advantages of GDA

- ② i.e. under what! when you do $P(x|y)$, you get log loss!
- but with a different step function



gives a new piece of evidence when can (learned) $\rightarrow P(x|y)$
 Having estimated all of those parameters, where

$$\text{MLE: } \hat{p}_y = \frac{1}{n} \sum_{i=1}^n I\{y_i=1\}$$

in terms of spans equals total have the word
 Denominator: number of spans equals
 all times you saw word ! if $y=1$
 Numerator: now: can turn feature back for

$$I\{y_i=1\} = \sum_{j=1}^n I\{x_{ij}=1\}$$

$$\hat{p}_{(x_i|y)} = \frac{\sum_{j=1}^n I\{x_{ij}=1\}}{I\{y_i=1\}}$$

Given a feature set write just likelihood of

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

Joint likelihood:

$$P(y) \prod_{i=1}^n p(x_i | y_i)$$

$$P(y) = P(y_1) \cdot P(y_2) \cdots P(y_{50k})$$

$$P(x_1, x_2, \dots, x_{50k} | y) = P(x_1 | y) \cdot P(x_2 | y) \cdots \text{Chain rule}$$

This is a false assumption, because words are related.
 But Naive Bayes works well for text docs.

$$\hat{p}_{(x_i|y)} = \frac{1}{I\{y_i=1\}} \sum_{j=1}^n I\{x_{ij}=1\}$$

$$\text{Assume: } x_i \text{ are conditionally independent given } y$$

$$N = 50,000 \quad b/f \text{ each} \rightarrow 2,50,000 \text{ parameters for multi-variate dist'n}$$

$$y \in \{0, 1\}^n : binary value$$

$$p(x|y)$$

Let's build a generative model for this

symmetry - feature

buy	0
car	1
spouse	0
adult	1
symmetry	0

$x = [1, 0, 1, 0, 0, 1]^T$ feature vector

$y \in \{0, 1\}^n$ e.g. email spam

Naive Bayes

N_{95}

1 # + 1 # 15
 $\frac{1}{S_1 \#} - P(y=1)$
 MLG for 3/15 sequential random variable
 : Louisiana
 : U.S.
 2/22
 0
 0
 2/11 W.A.S.
 2/8 W.S.
 2/8 Bonferroni's lemma $\frac{1}{S_{\text{total}}}$
 Fix:

just : you have it now + before
 $P(x_{300} | y=1)$

still probably bad idea +
 $\frac{o+o}{o} = so$

undetermined \Rightarrow

$$o = (1=h) \stackrel{\text{so}}{=} P(x|y=1) + (1=h) \stackrel{\text{so}}{=} P(x|y=0)$$

$$P(y=1|x) = \frac{P(x|y=1)P(y=1)}{P(x|y=1)P(y=1) + P(x|y=0)P(y=0)}$$

so span changes

$$o = (o=h) \stackrel{\text{so}}{=} P(x_{300} | y=0)$$

$P(x_{300} | y=1) = o = (1=h) \stackrel{\text{so}}{=} y_{\text{new}}$, we never see it before!

e.g. NIPS conference

so it turns out that this idea always works

95

Laplace smoothing

$$\frac{1 + nr}{1 + \sum_{i=1}^n} = (1=h) \stackrel{\text{so}}{=} P(y=1) + h + \text{more stuff}$$

More stuff

$$\left[\frac{1+o+o+1}{1+o} = P(w^n) \right]$$

Add 1 to all counts

\therefore will be sum more terms?
 Will we sum more terms?
 More stuff

$$\left(\frac{1+nr}{1+\sum_{i=1}^n} \right) = (1=h) \stackrel{\text{so}}{=} P(y=1) + h + \text{more stuff}$$

(multinomial MLG)

(4)

you maximize this in terms of posteriors.

$$\log \prod_{i=1}^m P(y_i | \phi) = \log \prod_{i=1}^m P(x_i | \phi, k_i) = \log \prod_{i=1}^m P(x_i | \phi, k_i = 1) + P(k_i = 0 | \phi)$$

$$P(k_i = 1 | \phi, y_i = 0, x_i) = P(k_i = 0 | \phi, y_i)$$

write down the likelihood of posteriors

What is the def. of MLE?

that we're the word k_i

Ratio: word is the fraction of all words in span

- total length of all spans
- sum over fractions of spans, use the length

Denominator

all the times you decoded word k_i

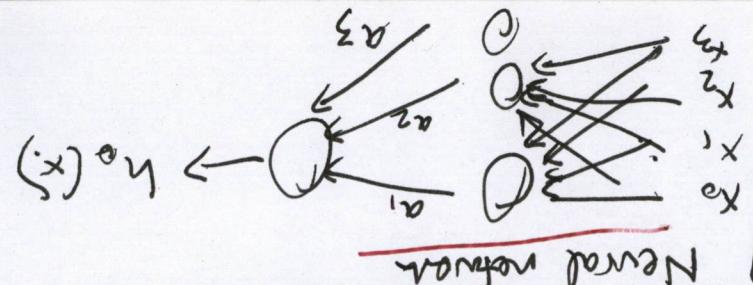
- sum over all your counts, & take open, x count of

Numerator: $\sum_{i=1}^m P(y_i | \phi, x_i)$

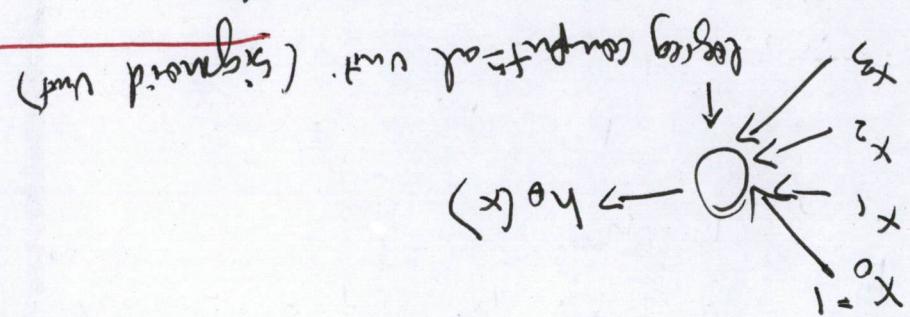
$$\frac{\sum_{i=1}^m P(y_i | \phi, x_i)}{\sum_{i=1}^m P(y_i | \phi, k_i = 1, x_i)}$$

return a training set, we can work out the MLE

(Nq 6.)



To get non-linear decision boundaries:



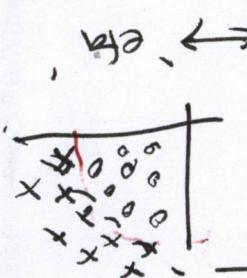
User NB, actually want a linear classifier.

$$x_i y_i \sim \text{exp}(\theta_i)$$

Is there a non-linear boundary?

$$\text{logistic } h_0(x) = \frac{1}{1 + e^{-\theta_i x_i}}$$

Non-linear classifiers



- For left class, the multivariate error of model does better: The next model takes into account the number of times it appears, but not sure why...
- The bigger & firmer don't do better as classifying

Video: "The machine that changes things"

The "Neuralink" neural network to read thoughts

Very "messy": generalised extreme

Scenarios: use child voice

"difficult" easily and quickly

Support Vector Machines (SVM)

After algo to get non-linear classifiers

$\text{start} = \text{linear! linear separator}$

• Predicts if $x_1 > 0$

• Counts if $x_1 < 0$

If $x_1 < 0$, then very confident that $y = 1$

$x_1 < x_2 \theta_1 + \theta_2 x_2 + b = 0$

Assuming linear separability

Def is most distance

"geometric margin"

Summation idea for digit recognition

You can look at word hidden layer is doing a.k.a. "back propagation"

Predicts if a_i

use grad descent to make it better

$J(\theta) = \frac{1}{2} \sum_{i=1}^{n_{\text{train}}} (y_i - h_{\theta}(x_i))^2$

$h_{\theta}(x) = g(x^T \theta)$ where $\theta = [a_0 \dots a_n]$

$a_i = g(x^T \theta^{(i)})$ where $g(z) = \frac{1}{1+e^{-z}}$

Each unit has its own parameter

$N_{\text{grad}} = 2$

For neural network, lots of local optima!

Out next researches support sum as better

of the day!

Several networks use the best for now

video year

Robustness of noise

3

Videos: "the machine that changes things"

the "world" BBC / PBS

to "Neuralink" neural network to read thoughts

Each unit has its own parameter

$N_{\text{grad}} = 2$

note diff notation



$$\left[\frac{\|w\|}{q} + \langle x_i, \frac{\|w\|}{\|w\|} \rangle \right]_{(i)} = y_{(i)}$$

more generally, define geometric margin

If you have training example $x_{(i)}$, then can compute the distance!

$$\frac{\|w\|}{q} + \langle x_{(i)}, \frac{\|w\|}{\|w\|} \rangle = y_{(i)} \quad \dots$$

$$\|w\|_1 \lambda = \frac{\|w\|}{\|w\| w^T \mathbf{1}} = q + \langle x_{(i)}, w \rangle$$

hyperplane

the separation for the

$$s.t. y_{(i)} \geq 1 - \epsilon \text{ s.t. soft margin}$$

the hyperplane = $\frac{\|w\|}{\|w\|} \cdot \lambda - \langle x, w \rangle$

unit length: $\frac{\|w\|}{\|w\|} = \|w\|^{-1}$

$$0 = q + \langle x, w \rangle$$

geometric margin

\therefore Need a normalization condition? $\|w\| = 1$

just say $w \leftarrow \frac{1}{\|w\|} w; b \leftarrow \frac{b}{\|w\|}$.

To make function margin large, could

- we want it to be large

margin is the margin

as far as possible

if $y_{(i)} (w^T x_{(i)} + b) > 0$, then correctly

$Q \gg q + \langle x_{(i)}, w \rangle$ and $w^T x_{(i)} < 0$

$Q \ll q + \langle x_{(i)}, w \rangle$ and $w^T x_{(i)} > 0$

Defn: functional margin of hyperplane

$w \in \mathbb{R}^{n+1}$

$x \in \mathbb{R}^n$ as input

e.g. separate out b

$g(x) = g(x^T w + b)$

$g(x) = \begin{cases} 1 & \text{if } x^T w + b \geq 0 \\ -1 & \text{if } x^T w + b < 0 \end{cases}$

have 4 output values in

$\{ -1, 0, 1 \}$

change notation

Ng 6

of SVM !
 Will directly info infinite dimensional space
 runs as well as (e.g.)
 shows quadratic margin
 You can scale w.b. by any factor, if it doesn't

$$\max_l \quad l \quad \text{such that } y_{(i)} (w^T x_{(i)} + b) \geq 1 = \|w\|_2$$

algorithmic part max quadratic margin

Maximum margin classifier: (precursor of SVM)

$$l = \min_i p_{(i)}$$

(3) quadratic margin

$$l = \frac{\|w\|_2}{p_{(i)}}$$

$$l = \frac{1}{\|w\|_2} \quad \text{if } \quad \text{facts}$$

$$N \in \sqrt{6}$$

- Chase w, b to max $\max_w \parallel w \parallel$: #1

$$\#1 \quad \max_w \parallel w \parallel \text{ s.t. } y_{(i)}^T (w^T x_{(i)} + b) \geq 1$$

$$w = \frac{\sum_i y_{(i)} x_{(i)}}{\sum_i y_{(i)}}$$

- But this is nasty : $\parallel w \parallel = \sqrt{\sum_i w_i^2}$
 This is a sphere! This is a non-convex
 constraint! Local minimum, etc

- difficult way of passing the same problem
 - hard as "find" from left $\{y\}$
 - fun = union of every example $y_{(i)}$
 - still mostly given margin
Suppose I have sufficient for convex optimization
it will classify to be big or possible, how to
by constraints (i.e. the support function)
new scaling ! $y_{(i)} (w^T x_{(i)} + b) \geq 1$

- SVM dual problem ! problem (KKT)
- optimal margin classifier
- total margin $M = \sum_{i=1}^n y_i$
- SVM dual and "margin optimization"
- kernels
- Naive Bayes

longer face will help to
keep out

Otherwise, $\Theta^*(w) = f(w)$: it's ~~satisfy~~ consistent!

$\nexists n: \psi(n) \neq 0$, then $\psi(p(n)) = 0$

$$\text{Max } \Theta^P(m) = \max_{\theta} \Theta^P(\theta) \quad \text{if } g(\theta) \leq 0 \quad \text{and} \\ \text{Min } \Theta^P(m) = \min_{\theta} \Theta^P(\theta) \quad \text{if } g(\theta) > 0$$

~~purple pants~~

$$\text{Logarithm: } L(a, b, c) = \sum_{k=1}^{\infty} \frac{(-1)^{k+1}}{k} (a^k - b^k + c^k)$$

$$\begin{aligned} \left(\begin{array}{l} \text{...} \\ \text{...} \end{array} \right) &= (m) \mathbf{u}_{\text{...}} & \mathbf{f}^{\text{...}} \mathbf{1} = ! & \mathbf{0} = (n)^{\text{!}} \mathbf{y} \\ \left(\begin{array}{l} \text{...} \\ \text{...} \end{array} \right) &= (m) \mathbf{b}_{\text{...}} & \mathbf{x} \cdots \mathbf{1} = ! & \therefore \mathbf{0} \leq (n)^{\text{!}} \mathbf{b} + \mathbf{g} \mathbf{n} \mathbf{s} \\ &&& \underbrace{\mathbf{(m) f u n c t i o n s)}_{\text{...}}} \end{aligned}$$

We form a general rule to solve a more difficult

$$\textcircled{2} \quad 0 = \frac{8C}{(8+m)8C} \quad ; \quad 0 = \frac{mC}{(8+m)8C}$$

for m to be a 5's necessary part
 $0 = \frac{\partial e}{\partial c} : 0 = \frac{\partial e}{\partial c}$
 long division
 $= 8$ $\overbrace{(m): 183}^{183 \neq (m)} + (n) = \cancel{(8)(m)}$

$$0 = \begin{bmatrix} (m)Y \\ \vdots \\ (m)Y \end{bmatrix} = (m)Y \quad \dots = ?'0 = (m)'Y + s$$

Primal & dual optimality

Take this, & throw off the silver
no local option! really sellable
i.e. a current problem faced now

Linear constraints had grown to

confirms
the notion
of a
paradigm

$$1 \leq (g + (\beta x_1 m))_{(1)} h = f \cdot s$$

$\Rightarrow \|M\|_2 \leq \min_{m \in M} \|f\|_2$

Add this to #2

t^{bN}

Kazishki kuru - Tüür

$$\text{Ansatz: } \alpha = f(w^*, x^*, \beta) = 0$$

$\frac{\partial}{\partial w} f(w^*, x^*, \beta) = 0$

$\frac{\partial}{\partial x} f(w^*, x^*, \beta) = 0$

$\frac{\partial}{\partial \beta} f(w^*, x^*, \beta) = 0$

with $x^* = g(w^*, \alpha^*)$

\downarrow

$\alpha = f(w^*, x^*, \beta)$

Then $E \bar{w} * \alpha + \beta + s.t. w * \text{sales}$ are parallel problems and depend on each other.

and suppose g : are strictly feasible

Conditions for d^* = $\frac{\text{Hessian } H \leq 0}{\text{Let } f \text{ be convex}}$ Suppose h_i is affine $[h_i(w) = a_i^T w + b]$

While down period passing, while down should all people
have down days also

! if $y = x$ then $\max\{x, y\} = x$
 ! if $y \neq x$ then $\max\{x, y\} = y$

Margit Lamp

d

$$(x^{\alpha})^{\beta} \theta^{0\leq \alpha < \beta}$$

—
—
—

卷之三

卷之二

三九四

• 118 •

(pcm)

70

$$0 = \sum_{i=1}^n \alpha_i y_i x_i$$

Dual problem

$$\begin{aligned} W(\alpha) &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i y_j \alpha_i \alpha_j \langle x_i, x_j \rangle \\ &= \underbrace{\langle w^T x, x \rangle}_{\text{inner product}} - \frac{1}{2} \sum_{i=1}^n \alpha_i^2 \underbrace{\langle w^T x_i, x_i \rangle}_{\text{norm}} \\ &= \left(\langle w^T x_i, x_i \rangle \right) - \frac{1}{2} \sum_{i=1}^n \alpha_i^2 \\ &= \frac{1}{2} \|w\|^2 - \underbrace{\sum_{i=1}^n \alpha_i}_{\text{Dual variable}} \# \end{aligned}$$

Take first 2 constraints, figure out plus

$$0 = \sum_{i=1}^n \alpha_i y_i \Leftrightarrow \alpha = 0$$

$$\boxed{\sum_{i=1}^n \alpha_i y_i x_i = 0} \Leftrightarrow$$

$$0 = \sum_{i=1}^n \alpha_i y_i - b = f(w)$$

$$\theta = \min_w \frac{f(w)}{\|w\|^2}$$

Dual problem:

$$(1 - (q + \alpha^T w)) \cdot \sum_{i=1}^n \alpha_i y_i - \frac{1}{2} \|w\|^2 = f(w) \#$$

$\alpha_i = 0$ for non-support vectors
There are few of them!

These examples are the support vectors
examples: fundamental margin $\gamma = 1$ if $\alpha_i \neq 0$

$$0 = q + \sum_i \alpha_i x_i \quad \begin{array}{c} x_1 \\ x_2 \\ \vdots \\ x_n \end{array}$$

$$T = \min_w f(w) \Leftrightarrow$$

$\alpha_i \geq 0$

$$0 = (q, w) : b \Leftrightarrow$$

$$0 \geq 1 + (q + \alpha^T w) \Leftrightarrow b = (q, w) : b$$

$$\dots \Rightarrow 1 \leq (q + \alpha^T w) \Leftrightarrow \frac{s.t.}{\|w\|^2} \leq \frac{1}{2} \|w\|^2$$

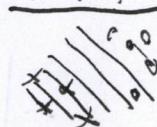
parameters: $w : \begin{cases} q, b \\ \alpha_i \end{cases}$

Logistic multiclass: $\alpha_i \in \mathbb{R}$

Applying of optimal design classifier of SVM

$$f(N)$$

- use of these inner products
 - you can make prediction by inner product
 - $\langle \vec{x}_i, \vec{w} \rangle$
 - kernels allow efficient computation if inner product
 - This holds only for certain feature spaces
 - like these multiclassify
 $\langle \vec{x}_i, \vec{x}_j \rangle$
 that will allow you to compute inner products
 there will be an implicit representation
 But despite that this it will turn out that
 features may be very much different
Kernels:
 various you are predicting
 $y = \sum_{i=1}^n \alpha_i y_i \vec{x}_i$
 $y = \sum_{i=1}^n \alpha_i y_i \vec{x}_i^T \vec{w}$
 Express surface larger as linear problems
 $\vec{w} = \sum_{i=1}^n \alpha_i y_i \vec{x}_i$
For next lecture

All you
 have only for support vectors. w will
 be perpendicular only along support vectors
 $\boxed{\text{Q3: } \alpha_i \neq 0 \text{ only for support vectors. } w \text{ will}}$
~~you can do anything to separate x_1 & x_2 without without~~
~~any condition~~
 $\boxed{\text{Q2: } w(x) : \text{you can do anything to separate } x_1 \text{ & } x_2}$
 $w = \max_{i=1}^4 w_i x_i + \min_{i=1}^4 w_i x_i$
 i.e. find max
 such that
 + best in either
 $w = \sum_{i=1}^4 \alpha_i y_i x_i$

 for b :
 Due to max for α_i , we get b ! than easy to see
 for SVM, we'll take dual of previous for α
 So, we developed a dual optimization problem
 $\boxed{\text{Q4: } \max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \vec{x}_i^T \vec{x}_j}$
 $\text{such that } \sum_{i=1}^n \alpha_i y_i = 0, \text{ then } \theta^D(\alpha) = \sum_{i=1}^n \alpha_i y_i \vec{x}_i$
 $\infty = (\alpha) \theta^D \left(\sum_{i=1}^n \alpha_i y_i \vec{x}_i \right)$
 $\text{if } \sum_{i=1}^n \alpha_i y_i \neq 0$
 t, b

Idea of SVM is to find separation plane in feature space $\phi(x)$, such that $\phi(x)$ is high dim.

feature ϕ finds lots of words in feature space (inner products), result is a kernel

of compute it very inexpensively

$$\langle \phi(x_i), \phi(x_j) \rangle = K(x_i, x_j)$$

we can write down a kernel function

"If terms off in many important special cases" "If dimensions are large you can't compute them"

so dimensions. Then you can't compute them;

so features $\phi(x)$ is very high dimensional, sometimes

Repeating $\langle x_i, x_j \rangle$ with $\langle \phi(x_i), \phi(x_j) \rangle$

$$(x)\phi = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} \leftarrow x$$

Take $x \in \mathbb{R}$ if to

have $x \in \mathbb{R}$ e.g. linear area

The only dependence of algorithm on x is from the inner products (of two input features)

Kernels

$$g(\sum x_i y_i x_i^T + b)$$

$$(w^T x + b)$$

inner prod

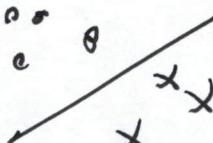


$$\max_{\alpha} \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j x_i^T x_j$$

Dual problem

margin

classifies the



$$s + y_i (w^T x_i + b)$$

Recap

- some algorithms

- soft margins (non-linear)

SVM - kernels

SVM - cuttines

Ng 8] $\text{if } L \leq \text{views}$

$K(x, z) = \log(1 + e^{-\|x-z\|})$ - "soft" -
which is smooth

(Centrifugal) one way is to try to come up
of If I gave you some random things
inner product will be large!

so if x, z are similar, then $\phi(x), \phi(z)$ will
be similarly in the same direction! so if

$$\langle \phi(x), \phi(z) \rangle$$

$$\frac{\langle \phi(x), \phi(z) \rangle}{\text{inertia: distance}}$$

How do you come up with a kernel?

High dimensional

(inner products, & then implicitly works)
using dimension feature vector! just compute

up to degree d

$(n+d)$ features of all monomials

$$K(x, z) = (x^T z + c)$$

other c.g. of kernels

②

$$\text{Define } K(x, z) = (x^T z + c)^2$$

$$\begin{bmatrix} x_1 & x_2 \\ x_2 & x_3 \\ x_3 & x_1 \\ x_1 & x_2 \\ x_2 & x_3 \\ x_3 & x_1 \\ x_1 & x_3 \\ x_3 & x_2 \\ x_2 & x_1 \end{bmatrix} = \phi(x)$$

compute $K(x, z)$
but instead $O(n^2)$ time to

$\phi(x)$
Need $O(n^2)$ to compute

$K(x, z)$ corresponds to a feature mapping

$$((z)\phi)((\phi(x))) = (\sum_{i=1}^d z_i x_i)(\sum_{j=1}^d x_j z_j) =$$

$$(\sum_{i=1}^d z_i x_i)(\sum_{j=1}^d x_j z_j) = (x^T z)^2 = K(x, z)$$

$x, z \in \mathbb{R}^n$
c.g. 2 inputs

kernel

Ng N

i.e. positive semi-definite matrix

$$= \sum_{i,j} k(x_i, x_j) z_i z_j \geq 0$$

$$= \sum_{i,j} \sum_{i,j} z_i (\phi(x_i))^\top (\phi(x_j)) z_j$$

$$= \sum_{i,j} (\phi(x_i))^\top (\phi(x_j)) \sum_{i,j} z_i z_j$$

$$= \sum_{i,j} z_i \phi(x_i)^\top \phi(x_j) z_j$$

$$z^\top K z = \sum_{i,j} z_i k(x_i, z_j)$$

Then, for any vector $z \in \mathbb{R}^n$

$$k(z) = k(x_0, z)$$

be given: Let $K \in \mathbb{R}^{n \times n}$ matrix

Now consider this is a valid kernel, let $\{x_0, \dots, x_n\}$

then k is valid matrix $K \in \mathbb{R}^{n \times n}$

if for all $\{x_0, \dots, x_n\}, n > \infty$

exist some ϕ s.t. $k(x_i, z) = \langle \phi(x_i), \phi(z) \rangle$

Is this really a valid kernel? Does there exist

$$k(x, z) = \exp(-\frac{\|x - z\|^2}{2\sigma^2})$$

Convex also holds

Theorem (Hilbert). Let $k(x, z) \in \mathbb{R}$. Then

$$k(x, z) = \phi(x)^\top \phi(z)$$

k is a valid (inner) product i.e. $\in \mathcal{P}$

Converse also holds

$$\text{iff } k(x_i, z) = \phi(x_i)^\top \phi(z)$$

if for all $\{x_0, \dots, x_n\}, n > \infty$

then k is valid matrix $K \in \mathbb{R}^{n \times n}$

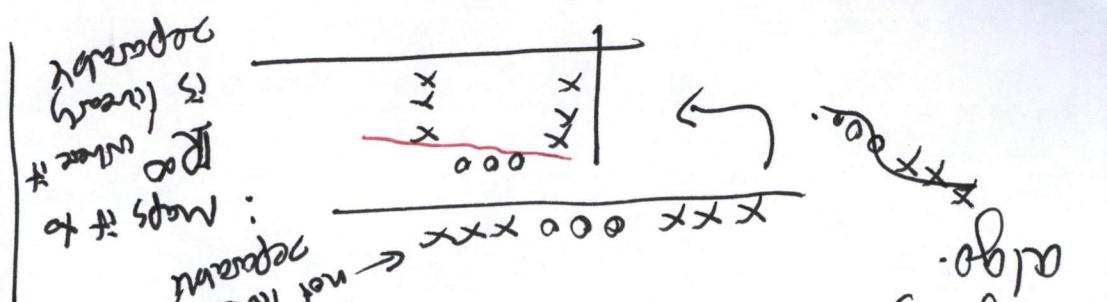
$$k(x_i, z) = \phi(x_i)^\top \phi(z)$$

k is a valid (inner) product i.e. $\in \mathcal{P}$

Theorem (Hilbert). Let $k(x, z) \in \mathbb{R}$. Then

(Kernel trick)
 the algorithm is more in the dual form
 can relate with $K(x_i, x_j)$, & how
 in terms of those inner products. Then you
 can take most algorithms of permute from
 algorithm using <> products. then
 the dual & allowed us to write the entire
 idea of kernels is more general than SVM.

not separable.
 separable. This algorithm won't work if
 you do + know about them are linear
 to do is solve convex optimisation problems
 boundaries. In the entire process, all you need
 This is how they output now. linear decision
 regions (hyperplane) to the largest possible margin
 machine (classifier). In the classifier and
 SVM is as oppose, find the optimal
 (4)

To the basic of SVM
 e.g choose $K(x, z) = \exp\left(-\frac{\|x-z\|^2}{2\sigma^2}\right)$ (Gaussian kernel)
 or $(x^T z + c)^d$, etc
 close depends on the problem.
 go back to dual form:
 Replicat $\langle x_i, x_j \rangle$ with $K(x_i, x_j)$
 Replat feature vector + with high dim
 Thus a Gaussian kernel is an as dual
 feature vector.
 future vector. If you can run sum in
 future time. You never need of compute
 of dual feature vectors. Why is this a good idea? We stated out
 why we would never linear learning
 algorithm.


These are convex/concave conditions

$$1 = \{y_i(w^T x_i + b) \leq 0\} \iff 0 < \alpha_i \leq C$$

$$\alpha_i = C \iff y_i(w^T x_i + b) = 1$$

$$\alpha_i = 0 \iff y_i(w^T x_i + b) \leq 1$$

What constraint is optimum?

Use KKT conditions! How do we know

Ansatz

$$\dots = ? \quad \alpha_i \geq 0$$

$$0 = \frac{\partial}{\partial \alpha} L(\alpha)$$

$$\max_w w^T x = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \alpha_i (y_i w^T y_i)$$

$$-\sum_{i=1}^n \alpha_i y_i = 0$$

$$\left(\sum_{i=1}^n \alpha_i - 1 - (w^T x_0 + b) \right) = 0$$

$$\mathcal{L}(w, b, \alpha, r) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \alpha_i$$

Dual of this problem can be done

If this isn't fast

②

time, but $C \gg \dots$ penalizes for large $\|\alpha\|$

∴ Allowing SVM to misclassify some of the constraints correctly!

Penalty: if $y_i(w^T x_i + b) > 0$,
by separating by $1 - \varepsilon_i$.

terms ε_i ? if demand that each training example

modifies by adding ε_i in red! The penalty

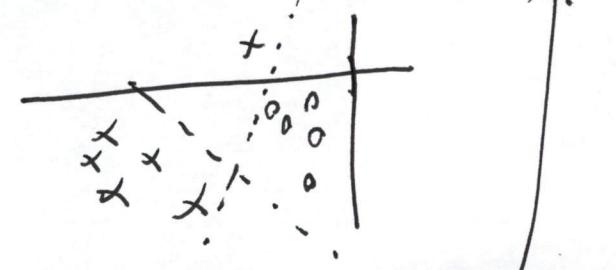
$$0 < \varepsilon_i \leq 1 - \sum_{j=1}^n \alpha_j y_j$$

$$\min_w \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \varepsilon_i$$

SVM primal problem uses:

is not linearly sep.

deals = the bottom x! or with dots for



L1 non soft margin SVM

Non-linear soft margin SVM

[N98]

⑥ Optimize $W(\alpha_1, \alpha_2, \dots, \alpha_m)$ ref. to α
Apply to Sym and optimsn: SMO
 in boost from CA direction
 $\sum_{i=1}^m y_i \alpha_i = 0$; if you fix all α_i 's but one, you can't
Sequential minimal optimisation: (SMO)
- Choose 2 α_i at a time
Outline
 Hold everything except α_i fixed
 Hold α_i fixed except α_j
 {optimise $W(\alpha)$ w.r.t α_i , s.t. $\sum_i y_i \alpha_i = 0$
 Run until you satisfy convergence
 constaint
 Iter loop can be done efficiently

May have heuristic to determine what else to update α :
 (updated to N. Witten's advice, CA w.r.t first
 max step. But different interpretation of
 max step)



Hold everything except α_i fixed:
 $\alpha_i := \arg \max_{\alpha_i} W(\alpha_1, \dots, \alpha_i, \dots, \alpha_m)$
 for $i \in 1 \dots m$
 repeat {
 to converge
 the problem $W(\alpha_1, \dots, \alpha_m)$, no constraints
 }

Problem : coordinate descent
Algorithm for solving optim.
 Ng 8.6

Cout of calculations of softmax :

$$\text{softmax} [x] = \frac{e^x}{\sum e^x}$$

Want down all \rightarrow count of 4 down acids

$$d(x) = ?$$

Assymetric function to down

reflect: BAT STAB AT BAT

e.g. protein acc. classification
20 amino acids, A ... Z

The pixels could be different.

But SVM doesn't know about acidity,

$$= \alpha_i (-\|x - z_i\|^2)$$

constant

$$K(x_i, y) = (x_i^\top y)$$

Σ

Normal network was complete

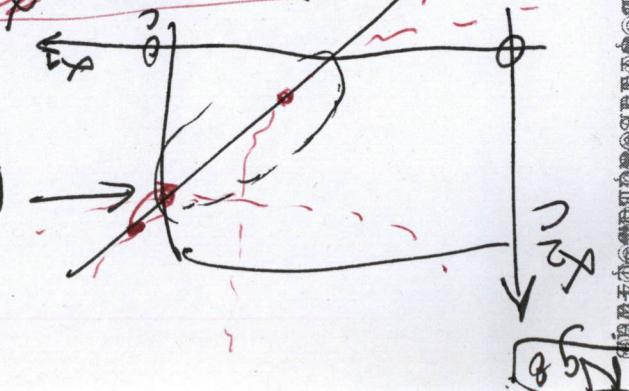
$x \in \mathbb{R}^{100}$: 100 features



e.g. handwritten digit recognition

To wrap up: applications

Dot constraints



Efficient for classification $\propto \alpha_2$ is computation

Easy to calculate!

On every iteration, choose $z \alpha_2$ to satisfy the source!

A D quadratic

$$= \alpha_2^2 + b\alpha_2 + c \quad ! \text{ easy of software}$$

$$W(\alpha_1, \dots, \alpha_n) = \frac{y_{(1)}}{\alpha_1}, \frac{y_{(2)}}{\alpha_2}, \dots, \frac{y_{(n)}}{\alpha_n}$$

$$\alpha_i = \frac{y_{(i)}}{\alpha_2 + y_{(i)}}$$

SVM due to the non-linearity
the dual learning algorithm

Dual form: $\phi(x)^T \phi(z)$
a quadratic can
compute

$\in \mathbb{R}^{160,000}$ features

$\phi(x) \in \mathbb{R}^{(20,4)}$

Ng8

some diff classes of f

$$(y) \in \mathcal{E}^s(h)$$

As you vary θ , you get different f

from X to $\mathcal{D}_0, \mathcal{D}_1$

$$\text{Hypothesis } h_\theta = \theta_0 + \theta_1 x_1$$

see all along as "classifying from a class"

cover all classes.

Now - count of features = feature ! they find a way ! can be visualized as approximation of this log reg! & SVM are easy can be visualized thus

$$\theta = \arg \min_{\theta} \sum_{i=1}^m \ell(y_i, h_{\theta}(x_i))$$

equivalent to minimize

ERM algorithm

(a.k.a. ridge)

Training error of h_θ :

parametric function $g(x)$ $\left(\begin{array}{c} x_0 \\ x_1 \end{array} \right)$ \rightarrow 1D fit from some data

$$g(x) = \sum_{i=0}^n \theta_i x^i$$

e.g. linear classification equation

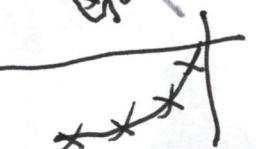
forward model of ML



classif. problems

for θ

logit model



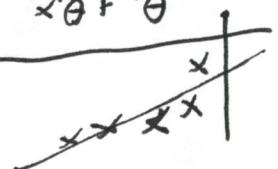
the data shows various patterns in

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_d x_d : \text{Overfitting} = \text{Wiggly Variable}$$

still doesn't fit.

$$\theta_0 + \theta_1 x_1 = \text{bias} : \text{Bias of data, but}$$

0.5 - wiggly bias



Bias - Variance tradeoff

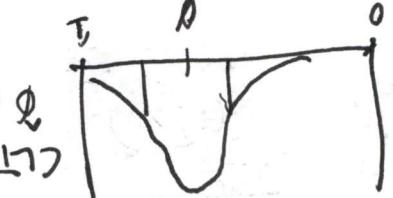
what quadratic people who really got ML... vs the person who need a feature

29.2.25

"Learning theory"

No. Lecture 19

As given in figure in the middle of file will
determine expansion



will be gaussian

Then $P(|\hat{y} - y| < \epsilon) \leq 2e^{-2\epsilon^2/m}$
Let $\phi = \frac{1}{\sqrt{m}}(z - \bar{z})$: if $|\phi| < 1$ then $\hat{y} = (1/\sqrt{m})z + \bar{y}$
Bernoulli (ϕ) random var $P(\hat{y} = i) = \phi^i$
Let $z_1, z_2, \dots, z_m \sim \text{iid}$ $E[\hat{y}] = \bar{y}$



$P(A_1 \cup A_2 \cup \dots \cup A_k) \geq P(A_1) + P(A_2) + \dots + P(A_k)$
(not necessarily independent) true
Union Bound: if A_1, A_2, \dots, A_k be k events,
then $P(A_1 \cup A_2 \cup \dots \cup A_k) \leq \sum_{i=1}^k P(A_i)$

Generally we care about the case where:

$$E(h) = P(h \sim H \mid (h(x) \neq y))$$

Generally: not training error, but pred

b99

for large m , will be bounded
Prone that for fixed y_i , h_i training error for the
 \hat{y}_i by definition is ineqvally:

$$\Pr[h_i \neq y_i] \leq \Pr[\hat{y}_i \neq y_i] = \Pr[\hat{y}_i \neq E[\hat{y}_i]]$$

which is

$$\Pr[\hat{y}_i \neq y_i] = \Pr[h_i \in \mathcal{H}]$$

$\Pr[h_i \in \mathcal{H}] \approx \Pr[h_i \in \mathcal{H}]$
if \mathcal{H} is finite

Want to show that ERM helps on generalization error

$$h = \arg \min_{h \in \mathcal{H}}$$

$$\mathcal{H} = \{h_1, h_2, \dots, h_n\}$$

The case of finite \mathcal{H}

②

Highly Probable

"Uniform convergence result"! as $m \downarrow, \mathbb{E}(h)$
will all converge to $\mathbb{E}(h)$ (almost surely)

will be uniform if $\mathbb{E}(h)$ for all $h \in \mathcal{H}$
so with prob $1 - 2e^{-2\gamma m}$ we have $\mathbb{E}(h) *$

$$\leq 1 - 2e^{-2\gamma m}$$

$$(P \leq |(h)\mathbb{E} - \mathbb{E}(h)|, \text{ if } h \in \mathcal{H}, \dots)$$

$P(\text{not } h \in \mathcal{H}) = 1 - \text{prob sides}$: i.e. these hypotheses = in case

$$\leq \sum_{i=1}^k 2 \exp(-2\gamma m) = -2k e^{-2\gamma m}$$

$$= P(A_1 \cup A_2 \cup \dots \cup A_k) \leq \sum_{i=1}^k P(A_i)$$

$P(E(h)) \leq \frac{1}{k} \sum_{i=1}^k P(A_i)$
(suppose that I make a loss
case)

$$P(A_i) \leq 2e^{-2\gamma m}$$

$$A_i = \{h_i \text{ is true}\}$$

Now prove for all B_i, h_i

659

Restate uniform convergence bound
(given δ, γ , what is m ? (finding m)
(given δ, γ , what is m ? (finding m)

$$\delta \leq |(h)\mathbb{E} - \mathbb{E}(h)| \text{ for all } h \in \mathcal{H}$$

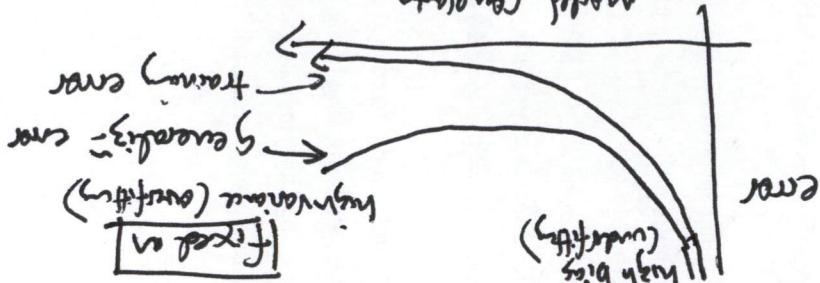
$$\text{so } (\delta/m) \leq m \leq \frac{1}{2\gamma} \log \frac{2}{\delta}$$

$$\delta = 2ke^{-2\gamma m}, \text{ solve for } m:$$

Restate uniform convergence bound

Can we prove $\mathbb{E}(\mathbb{E}(h)) = h$ = $\mathbb{E}(h)$
(unif. convergence)
Let's assume $\mathbb{E}(\mathbb{E}(h)) \leq h$
Last step of overall proof.
 $\frac{\mathbb{E}(\mathbb{E}(h)) - h}{\sqrt{2\gamma m}} > |(h)\mathbb{E} - \mathbb{E}(h)|$
So for f , for fixed m, g
"error bound":
as $\delta \rightarrow 0$
i.e. number of training examples do not grow fast

Consider the surface of quadratic, traffic.
 Let's do it have
 H
 f
 lower \rightarrow quadratic
 (we're)
 In eqn (4), left term is bias, 2nd is variance
 So coupling class will make variance zero.
 (for confidence problems, no square fully accepted)
 def of bias of variance - n.b
 $E(H)$ will become better ($b_0 + q_u$ adaptive)
 but k will increase (i.e. rise of net fitting is
 In eqn (4), left term is bias, 2nd is variance
 (so coupling class will make variance zero)



Left terms deal with bias - trade off ear
 which implies eqn (*)

$$l = \int_{\frac{1}{2} \ln \log \frac{g}{g}}^{\frac{1}{2} \ln \log \frac{g}{g}} \text{W.E. term eqn (1)} \text{ walls up-l-g}$$

$$\begin{aligned} & (*) \exists \\ & \left(\frac{1}{2} \ln \log \frac{g}{g} \right) \end{aligned}$$

Then w.p. 1-g : $l = |H| + f$: linear

$$\begin{aligned} & \text{to ears of } \exists = (\gamma) \exists \\ & \text{w.r.t. } \exists = (\gamma) \exists \\ & \text{w.r.t. } \exists = (\gamma) \exists \\ & g \cdot u \end{aligned}$$

- (1) $l = f + b + (*) \exists \Rightarrow$
- (2) $l = f + (\gamma) \exists \Rightarrow$
- (3) $l = f + (\gamma) \exists \Rightarrow (\gamma) \exists$

gen. ears:
 best hyperplanes for (3)

$$(\gamma) \exists = \arg \min_h$$

69

19

Corollary : Let $H = \{f\}$ be fixed
Then, suppose I want to guarantee
 $E(\hat{Y}) \leq \min E(Y) + \epsilon$
 n
w.p. $1 - \delta$. If sufficient that :
 $n \geq \frac{1}{\epsilon^2} \log \frac{2}{\delta}$
 $\# \text{continuous functions} \leq O\left(\frac{1}{\epsilon^2} \log \frac{1}{\delta}\right)$
of which follows from the fact
that this is a tree and
thus all

Let $|E| = k$, $\alpha f, g$ be fixed, then for $E \in \mathcal{B}(f \geq g)$, we have $E(f) \leq E(g)$

f is a linear function of all linear decision boundaries

e.g. if f is parallelized by real numbers

e.g. $f(x) = \frac{\text{sum of class } 1 - \text{sum of class } 2}{\text{sum of all classes}}$

and based on the number of training examples needed

n : training data size

k : size of hypothesis class

$E(n) = \text{generalization error}$

ERM: emp. risk. min : chooses h that min generalization error

that $n \geq \frac{2}{\epsilon^2} \log \frac{2}{\delta}$ $= O\left(\frac{1}{\epsilon^2} \log \frac{2}{\delta}\right)$

Designation standards & regulations

- feature selection
- D105-N1A - Cross-validation
- model selection -
- VC dimension

\therefore All the theory of ERM gives good information for what types of SVM are doing!

Group 6 of own this step first
but this is our convex
line (eg) SVM can be used on convex after
but this is our convex
"if this is not linear: using - loss fn"
SVM is perceivable later: using - loss fn

$$0 \simeq ((x)^y)_{\perp} \quad x_{\perp} \theta = z$$

$$((x_{\perp} \theta)G = (x)^{\theta} y$$

+): $\{ h \neq (x)^{\theta} y \} \xrightarrow{\exists}$

The ERM benefit of algorithms:

By now you know, if for small VC dimension

$$VC(X) \leq \lceil \frac{d+1}{2} \rceil + 1$$

If of large margin separator

~~$$\vdots \vdots \vdots \vdots$$~~

Is $VC = R$?

② Laying eggs would be negatively linear such
VC dimension! Very closely similar to numbers
of parameters! upper-bounded by VC dim;

$(p)^{\text{pf}} \circ = m$

To guarantee that if suffices that

Corollary

$$\left(\frac{g}{1} \log \frac{m}{1} + \frac{g}{m} \log \frac{m}{g} \right) \alpha + (\beta - 1) \beta \geq (\gamma)^3$$

Thus, we have
from left side
 $\frac{g}{1} \log \frac{m}{1} + \frac{g}{m} \log \frac{m}{g} \geq (\gamma)^3 - (\gamma)^3$

~~Acknowledgments~~ Best known result in learning theory ($O(18 \cdot n)$)

$$1 + v = \begin{cases} p - v & v \\ 0 & \text{otherwise} \end{cases}$$

More generally, in addition, the VC dimension cannot

E.g. It is claim of lower claimant in 2D

of the largest of staffordshire.

Defn: $V - C$ discrepancy: $AC(3f)$ is the way

Model Selection (M.S.)

[Nag 10]

M.S. above automatically handle off by bias

of variance



\leftarrow model config (era)

$$\begin{aligned} \theta_0 + \theta_1 x & \\ \theta_0 + \theta_1 x + \theta_2 x^2 & \\ \theta_0 + \theta_1 x + \dots + \theta_n x^n & \end{aligned}$$

e.g. \rightarrow - boundary parameter \rightarrow MLE

How to choose these models? automatically select

$$M = \{M_1, M_2, M_3, \dots\}$$

choose C in SVM in $\frac{1}{2} \|w\|^2$

idea: train them all & look for training error min \rightarrow diff models \rightarrow no, enough!

- more expressive \rightarrow ! fits

\rightarrow leave-one-out CV

$k=m$: just one out each time!

- O is useful! use for scarce training S.

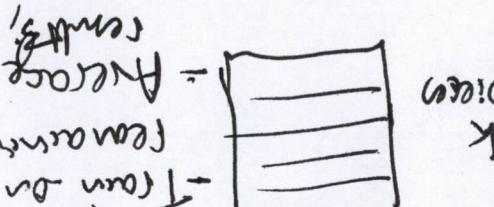
- more expressive O! fits

- k-fold CV

- more expressive O! fits

- $k=O \approx$ common! holdout 10%

- usually go back to train on whole S



- Train on $k-1$ pieces; test on the

- k-fold CV

- Accuracy over k-classes & their error

HCV works out great, but you often have very few training e.g. handful of held out N.

- optionally: retain this best model for the whole S.

- pick largest error in CV model.

- train each model on STrain, & test on

$S_{CV} (30\%) \rightarrow$ hold-out

S_{CV}

$\frac{\text{Hold-out Cross-Valid}}{\text{Bfits ideas}}$

$\frac{\text{Hold-out Cross-Valid}}{\text{S with S train (70\%)}}$

$\frac{\text{Bfits ideas}}{S_{CV} (30\%)}$

(3)

- NP-head paradigm of field-based feature test!

If you have 100 m² land then don't grow trees

all the features are off at first
 $f_1 = 13$ - $f_2 = ?$ steps -

Backward selection

cells forming a **raft**! conceptual expansions! efflux well;

Output best hypothesis found.

(3) For $i = 1 \dots n$, they add feature? if
for evaluate model using $CV = (\text{some folds})$

$$\emptyset = f_0 \text{ from parts -}$$

~~forward search | select : () feature of a file~~

Use secret websites of for foreign post)

Within 25 features, 2ⁿ possible subsets

Many ML have higher dimensional feature space, e.g. text classification ~ 50×10^3 p.
Reducing number of features might reduce variance, if right for overfitting
Select feature which tend to be relevant, & reduce risk of overfitting

F culture test tube

the people in England, the

- Taking of the loggers, it is in the process.

institu-

I find it best to have no place of storage on my compound so as to prevent 17th April

Learn: These boards are often used to hold
of wood for any kind of display! (including learning)
These boards will hang and display

132

• C15 NJ

Ng 10

how many features do we have
to choose k, we cross-validate to decide
how pick top k features - largest MI

$$MI(x, y) = \sum_{i=1}^n p(x_i, y) \log \frac{p(x_i, y)}{p(x_i)p(y)}$$

KL divergence! joint measure
of mutual information
of distribution of discrete
variables

$$MI(x, y) = \sum_{i=1}^n p(x_i, y) \log \frac{p(x_i, y)}{p(x_i)p(y)}$$

mutual information
of joint distribution
of discrete variables

e.g. corr(y, x_i)! pick top k most correlated features;
information Y? is about Y.
for each feature compute some measure of how

"filter" feature selection method
less computationally expensive:

⑤

① To predict a new point on y
 $P(y|x, s) = \int P(y|x, \theta)P(\theta|s)d\theta$
 Bayesian stats & regularization
 Unlike learning (dilution)
 Applying ML
 Reg: as a way of preventing overfitting
 i.e., take an expectation of y , given predictor
 If more predictors, this is computationally difficult.
 Instead of computing full posterior distribution
 $\theta_{\text{MLR}} = \arg \max_{\theta} P(\theta|s)$
 $(\theta) \left(\left(\theta \mid x_i \right) \mid P(y_i | x_i, \theta) \right) =$
 To make a prediction
 $\hat{y}_{\text{MLR}}(x) = \theta^{\text{MLR}} x$
 In practice, we will be close to θ . This is like Inertia: if the prior is $\theta \sim N(\theta^*, I)$, then instead matrix the prior?
 So, most won't be zero of size n : most likely to be close to θ^* . Thus: the feature selection, the most won't be zero of size n : most likely to be close to θ^* .
 vs Bayesian (equation)

We don't know θ :
 vs Bayesian school
 (i.e., frequentist procedure: ensure a true generating mechanism: θ is not random;
 calculate $P(s|\theta)$ for training set
 $s = \{x_i, y_i\}_m$
 $P(\theta) = \text{prior}$ E.g.
 $\theta \sim N(\theta^*, I)$
 $P(\theta|s) \propto P(s|\theta)P(\theta)$
 vs Bayesian school!
 Bayesian: calculate $P(s|\theta)$
 calculate $\max_{\theta} P(y_i|x_i, \theta)$
 e.g. linear regression
 choose θ such that likelihood
 i.e., frequentist procedure: ensure a true generating mechanism: θ is not random;
 vs Bayesian school

$$\sum_{i=1}^n \{y_i - \hat{y}_i\}^2 = \text{Total Sum of Squares}$$

Model with loops:
 All above so far are batch! without noise

Iteration: Online learning

(batch) Logistic regression becomes effective
 for text classifier, with boundary regularization

of own data to zero
 to fit surfaces, again \therefore max

$$w_i = |y_i - \theta^T x_i|^2 + \lambda \theta_i^2$$

With boundary (e)

$$\min_{\theta} \frac{1}{2} \|y - \theta^T x\|^2$$

Intuition: Minimize MLE

11/6/11

Add SCDS, or penultimate to do online learning
 Initialize $\theta = 0$
 after each update, run 1 step of SGD:
 $y_i(x_i) - \theta^T x_i$

$$\theta = \theta + \alpha (y_i - \theta^T x_i) x_i$$

You can prove that: even if

$X(i) \in \mathbb{R}^d$ (infinitesimal)

e.g., perpendicular will converge to perfect & define so long as some measure of separation is large enough

etc. the guys are in a no-discriminate space.

e.g., perpendicular will converge to perfect & define

as long as some measure of separation is large enough

for text classifier, with boundary regularization

of own data to zero

to fit surfaces, again \therefore max

With boundary (e)

$y_i - \theta^T x_i$

Intuition: Minimize MLE

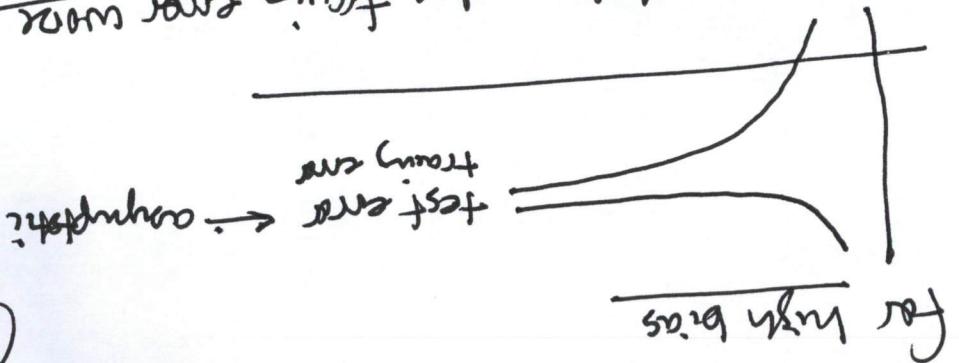
11/6/11

② Add SCDS, or penultimate to do online learning

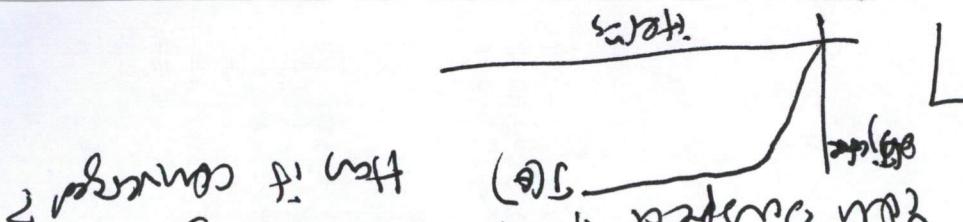
00:25 Applying ML algo
 I want to give you advice of how to apply really
 I mean (curve) of how much spreads out the
 points from the mean & wall. You just need
 to go through it well. (curve) of how much
 points fall on the wall. (curve) of how much
 of the data is far away from the mean &
 of the data is far away from the mean &
 of the data is far away from the mean &
 of the data is far away from the mean &
 of the data is far away from the mean &

3 key idea: - dimensionality
 - error of the last fitted model
 - how to select the latest sample

more training often makes training easier sooner
... try more exercises and more details, since that
can't help; This happens a lot. You could
have told them a month ago. You could
diagnose.
Dias vs varicous comes up often. A common
after need to consult your own diagnosis



B10s vs vacuums comes up often. A common
disadvantage
of the need to construct your own diagnostic
After exposure: SUM output from layers
Breaker log 59: 26 reseal for next cycle
Breaker log 62: 26 reseal for next cycle
10% of power! 0.1% of the
sum with 10% loss!



e.g. high degree of linkage: diagonal effect: linkage disequilibrium errors will be much lower than test errors will be much higher than linkage disequilibrium errors. Be aware of linkage disequilibrium: feature error will be smaller than test error: feature errors will be higher than test errors. For which variants, feature errors and gross with higher w2? So more training data will help, also smaller feature effects.

Better of figure and what of fix:
 But takes a lot time
 - my diff function for
 - get some terms off
 - sum off terms
 max $\sum_{i=1}^n \log p(y_i | x_i, \theta) - \lambda \|\theta\|^2$
 If you get 2% better test error
 Use log loss with grad descent

Debt-to-equity ML analysis

116N

- Case 2: $a(\theta_{SVM}) > a(\theta_{BLR})$
- If $J(\theta_{BLR}) > J(\theta_{SVM})$, then the problem is not linearly separable.
 - Let θ_{BLR} . If $J(\theta_{BLR}) < J(\theta_{SVM})$, then the BLR model is better than the SVM model.
 - If θ_{BLR} works in some parts but not all, then the SVM model is better.

Show it gives bad performance, what do you do?

$$\theta_{BLR} = \text{arg min}_{\theta} J(\theta)$$

use regularization term:

$$J(\theta) = ||x - f(\theta)||_2 + \lambda \theta^T \theta$$

Example of diagonalization for logistic

SVM ——— BLR

- for optimum algorithm = closed form
- more iterations = slow convergence

more to note! the weight vector
is orthogonal to most of the margins

\therefore BLR model $J(\theta)$, but SVM does better in

$$J(\theta_{SVM}) < J(\theta_{BLR})$$

$$\text{Case 2: } a(\theta_{SVM}) > a(\theta_{BLR})$$

(4)

diagonalization $J(\theta)$. \therefore BLR uses fewer iterations than SVM

But BLR uses many more $J(\theta)$. \therefore BLR

$$J(\theta_{SVM}) < J(\theta_{BLR})$$

Diagonalization $J(\theta_{SVM}) < J(\theta_{BLR})$?

$$J(\theta) = \sum_i \log p(y_i | x_i, \theta) - \lambda \theta^T \theta$$

BLR tries to maximize

$$a(\theta_{SVM}) < a(\theta_{BLR})$$

SVM outperforms BLR, so;

$$a(\theta) = \max \sum_i w_j h_j(x_i) = y_i$$

Weighed correctly: when class about:

BLR SVM outperforms BLR, but you need to

fix the regularizing? i.e. the regularization

Ask - what do you care about? Are your algorithms

No. 4

(5)

of tests for major sources of error

After first analysis
• this is often different b/w baseline & current performance.
• e.g. out-span counter, how lots of components + the reme counter & see how perform
• sharp

Gettings started on a learning path
#1: careful doors - best for reason
#2: build fix: a hole

- avoids preventive splitting
- faster to install.
- use early samples of diagnosis
- saves quality

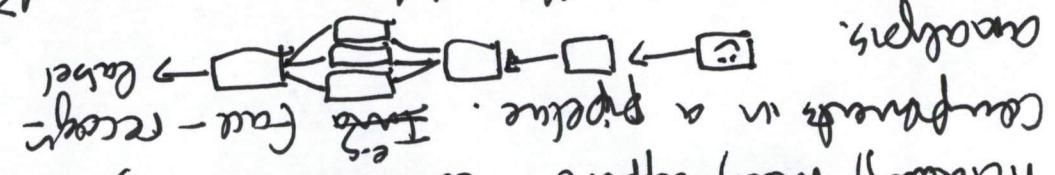
I use LR most often
find something wrong in the data
can think of: look at the future you
will do in every way you

Step 1: Do at the data every way you

use to figure out what is happening in an organization
some of you will graduate from
sharpie & get following high level applications
ML of sharpie: car & economy; apply
you should get an initiative understanding of
the problem! your own personal understanding
some silicon valley outsource their ML. often a
familiar idea. You don't keep the knowledge
diagnoses of early samples help justify the
program & your samples

Error from thresholds

see how accuracy changes? more fails,
Plug in ground truth for each component,
How much error is affordable of each component
measures.



[NG12]

Ng 11

Dates of over-theenzy
1995

Publication

Learning of VC dimensions
very important
of predictors
NP learners
With theory if I can see a link!

00:22:11

E=2 or 4?

The true number is ambiguous



try a few different k's:

$\text{Cluster } k^*$

can occur because two diff. clusters. Also (and often)

$\therefore J(c, \mu)$ decreases monotonically. K-means

can show that k-means is convergent and

$J(c, \mu) = \sum_{i=1}^n \|x_i - \mu_i\|^2$

for this to converge

Shows: Minimizes sum of squared distances from data points to cluster centers

(i) $\mu_i := \sum_{j=1}^n c_{ij} = \sum_{j=1}^n c_{ij} \times c_{ij}$

(ii) $c_{ij} := \sum_{k=1}^n \frac{1}{\|x_k - \mu_i\|^2}$

2. Repeat until convergence:

$\mu_1, \dots, \mu_k \in \mathbb{R}^n$: vectors in \mathbb{R}^n

1. Initialize cluster centers
(no clean loads, y)!

K-means algorithm

Input $\{x^{(1)}, x^{(2)}, \dots, x^{(N)}\}$

①

Aug 12

Unsupervised learning (UL) 28, 406

In UL, you are given a dataset with no labels. We also know of cluster structure



Applicability

Biology - genes

Market research - customers

Newsgroups - topics

Image segmentation

Image segmentation

The ML also has to do clustering structure

- EM (Expectation-Maximization)

- Gaussian's (inequality)

- Mixture of Gaussians

- Clustering (k-means)

Imagine there's a latent variable z . And $x^{(i)}, z^{(i)}$ have random variables
 Gaussian mixture: a joint (middle, unlabelled)
 density estimation
 (2)

$P(x^{(i)}, z^{(i)}) = P(x^{(i)}|z^{(i)})P(z^{(i)})$
 by the chain rule
 $z^{(i)} \sim \text{Multimodal}(\phi)$: (several for
 Gaussians)

$(\phi) \geq 0, \int \phi = 1$

$x^{(i)} | z^{(i)} \sim N(\mu_i, \Sigma_i)$, but y is dependent on z
 (it's not same as $x^{(i)}$, but y is dependent on z)
 π we know $z^{(i)}$, then we MLE.
 $\sum_i \log p(x^{(i)}, z^{(i)}) = \sum_i \log \phi(z^{(i)}) + \sum_i \log p(x^{(i)}|z^{(i)})$
 $\hat{\phi}(z^{(i)}) = \frac{1}{n} \sum_{i=1}^n \delta(z^{(i)})$
 $\hat{\mu}_i = \frac{1}{n} \sum_{i=1}^n x^{(i)}$
 $\hat{\Sigma}_i = \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \hat{\mu}_i)(x^{(i)} - \hat{\mu}_i)^T$

But we don't know z . So guess & using
 Expectation + MLE: EM

This diff is unusual w/ a standard one
 used in ~~accidently~~
 $P(y) \neq \text{the fact to detect anomaly}$
 $\text{Label training set } \{x^{(1)}, \dots, x^{(n)}\}$, but
 red flag: i.e. Anomaly detection
 if $P(x)$ is very small, that would result a
 loose of typical distn, build model of $P(x)$, then
 just change of model after
 you want to estimate
 for accurate enough:
 Head \rightarrow
 Density estimation
 Ny 12

which human eat that point same for
 anomalies first generated this. But I don't know
 could be mixture of 2 N : i.e. 2 separate
 dots

MLE of ϕ : (74) probability of getting ~~at least one~~

Recall - [this + explain - what happened]: we know ~~that~~ the cards:

COVARIANCE matrix for different location set

$$S_i = \frac{\sum_{j=1}^n w_j^{(i)} (x_j^{(i)} - \mu_i)^2}{\sum_{j=1}^n w_j^{(i)}}$$

In Gantt: we know ~~that~~ the details:

(HDS) to diff.

$$\{ \quad \quad \quad (3m) \Sigma) = 1$$

$$EM - \frac{\text{Expected}}{\text{Actual}}$$

E - Step
Expectation

Let $w_{ij} = \phi(z_{ij}) \times |x_i|$, $\forall i, j$.

E-step { given values of π_{ij} }
update counts

$$\geq P(x_{i_1} | z_{i_1} = \emptyset) P(z_{i_1} = \emptyset)$$

$$(\overline{(-\zeta)} \not\models (\overline{(-\zeta)} \rightarrow x) \not\models$$

- Computer graphics

三

$$\overline{(1) \times (1)^M} = : 1^n :$$

$$(3) \sum_{k=1}^{\infty} \frac{1}{k^2} = \dots$$

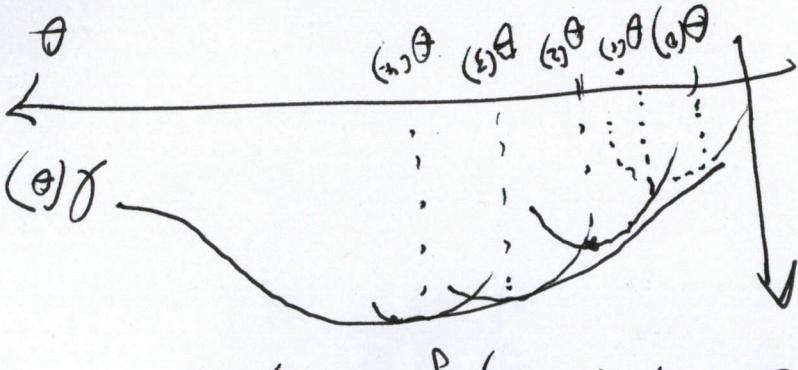
M-step: $\{ \text{update parameters} \}$

17

17

EM - EM

$z_1 b_N$



EM is a way of doing the ML above

i.e., averaging gradients and

$$Q(\theta) = \sum_{i=1}^n \log P(x_i, z_i | \theta)$$

$$Q(\theta) = \sum_{i=1}^n \log P(x_i, z_i | \theta)$$

Goal of maximization:

Differentiate only x .

Hence one useful for $P(x, z | \theta)$ - gradient

Derivative of gradient by

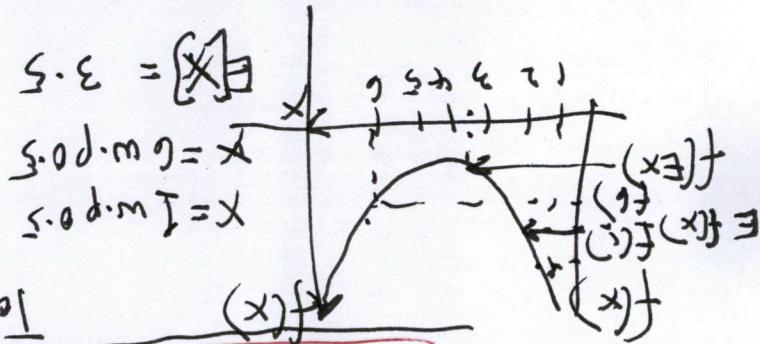
$$f(E[x]) \geq E[f(x)], \text{ etc.}$$

$$\text{If } f'(x) < 0 \text{ (convex), then}$$

thus $E[f] = f(E[x])$

Further, if $f''(x) > 0$, (f is strictly convex)

$E[x] = E[x]$ w.p. 1
 $\Leftrightarrow x$ is a constant w.p. 1



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

Jensen's Inequality

Let x " " a random variable

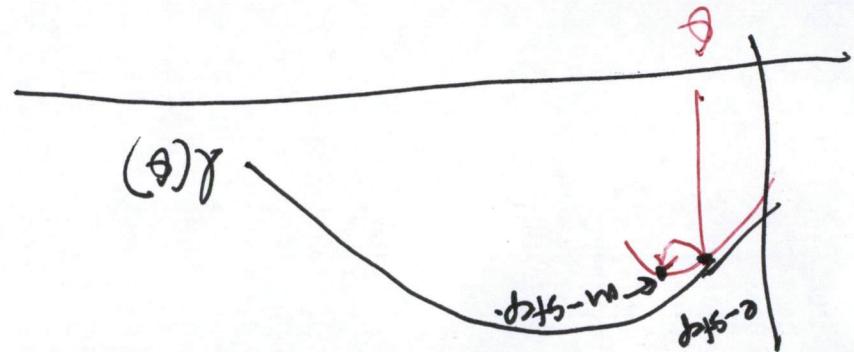
Means inequality

(Note, if N is a special case)

Gradual view of EM algorithm

21 N

Define what we do
 We constructed a lower bound for
 $\mathcal{L}(\theta)$. We show that:
 $\mathcal{L}(\theta) \geq \frac{\mathbb{E}[Z_{(1)}] - \mathbb{E}[Z_{(n)}]}{\mathbb{E}[Z_{(1)}] + \mathbb{E}[Z_{(n)}]} \cdot \mathbb{E}[Z_{(1)}]$
 The last two words is a
 bound
 inequality (as far as I can tell it's a tight bound)
 Now Jensen's inequality to be a
 - choose θ so that equality holds
 \Rightarrow for θ the second term is a constant, $\mathbb{E}[Z_{(1)}]$
 (remember $Z_{(1)}$ is a random variable)
 $\frac{\mathbb{E}[Z_{(1)}] - \mathbb{E}[Z_{(n)}]}{\mathbb{E}[Z_{(1)}] + \mathbb{E}[Z_{(n)}]} = \text{constant}$:
 for all values of $Z_{(1)}$
 $\mathbb{E}[Z_{(1)}] \leq P(x_{(1)}, z_{(1)}; \theta)$
 $\mathbb{E}[Z_{(1)}] =$



$$\frac{P(x_{(i)}, z_{(i)} | \theta)}{P(x_{(i)}, z_{(i)} | \theta)} = \frac{\prod_i P(z_{(i)} | \theta)}{\prod_i P(x_{(i)} | \theta)} = \text{constant}$$

- it's right lower bound

$$S \in \Theta : P(z_{(i)} | x_{(i)}; \theta)$$

$\hat{\theta}$:

EM algorithm iteration

$$P(z_{(i)} | x_{(i)}; \theta) =$$

$$\frac{P(x_{(i)}; \theta)}{P(x_{(i)}, z_{(i)}; \theta)}$$

$$= P(x_{(i)}, z_{(i)}; \theta)$$

$$\sum_{z_{(i)}} P(x_{(i)}, z_{(i)}; \theta)$$

see lecture
notes for
more details

$$\Theta : P(x_{(i)}, z_{(i)}; \theta)$$

Ng 12

⑥

$$\begin{aligned} & P(x_i | z_i) = P(z_i | x_i) P(x_i) \\ & \text{Use Bayes rule: } P(x_i | z_i) = \frac{P(z_i | x_i) P(x_i)}{P(z_i)} \\ & \text{EM alg: mixture of Gaussians} \\ & \text{E-step:} \\ & \quad P(z_i | x_i, \theta) = \dots \\ & \text{M-step:} \\ & \quad \text{EM is good around an } \hat{\theta} \end{aligned}$$

$\hat{\theta} \rightarrow \hat{z}(\theta, \alpha)$ (by Jaccard measure)

Define $J(\theta, \alpha) = \sum_i P(z_i | x_i, \theta, \alpha) \log \frac{P(x_i | z_i, \theta, \alpha)}{P(x_i | z_i)}$

Do this way of finding about $\hat{\theta}, \hat{\alpha}$

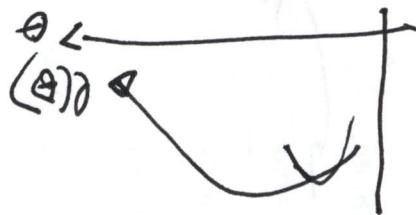
is tractable

e.g. for exponential family, $P(x_i | z_i, \theta)$

i.e. in closed form.

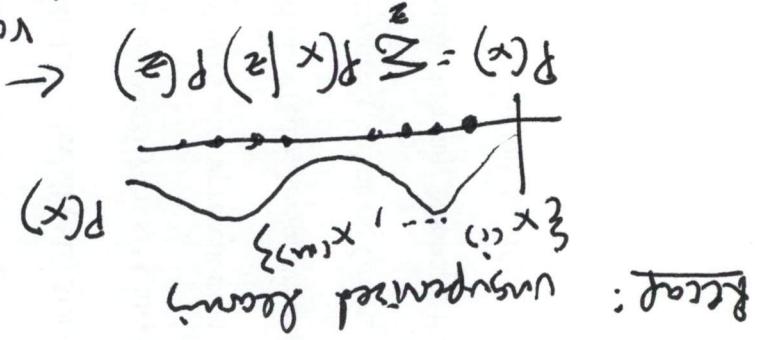
Because the M-step is tractable analytically

① Why do we settle to the EM algo?



$$\text{EM step: } \hat{\theta} = \arg \max_{\theta} \sum_i \log(P(x_i | z_i, \theta))$$

if came from
mixture of 2 Gaussians
versus flat distributions



Recap: nonparametric density

- Diagonalization: Gaussians dist

- factor analysis

- mixture of Gaussians

N 13

All in the lecture notes:

$$\theta = \frac{\partial \mathcal{L}}{\partial p}$$

Lagrange multipliers

$$(1 - \phi_i) + (\dots) = f$$

optimal

Lagrange

constraint equation for ϕ_i : construct the

Since ϕ is multinomial, $\phi_i = 1$ is a constraint.

do the updates

$$\Delta \mu_i(\dots) := 0 \quad | \quad \text{Solve for } \mu_i \quad | \quad \text{true for } \phi_i \rightarrow$$

Take derivative

$$\phi_i(z_{i,j})$$

$$w_{i,j} = \sum_{k=1}^K w_{i,k} \exp\left(\frac{z_{i,k}}{K}\right)$$

$$\log P(x_{i,j}, z_{i,j} | \theta, H_i)$$

$$\max_{\theta} \sum_{i=1}^n \sum_{j=1}^{K_i} \alpha_i(z_{i,j}) \log P(x_{i,j}, z_{i,j} | \theta, H_i)$$

$$M_{\text{step}}$$

Ng 13

Ex: feed clustering on in

news.google.com - group affected docs

apply EM:

Mixture of Naive Bayes

training set $\{x_{i,j}\} \in \mathbb{R}^{n \times K}$

feature of Naive Bayes

$x_{i,j} \in \{0, 1\}^2$ clusters

$$P(x_{i,j} | z_{i,j}) = \prod_{k=1}^2 P(x_{i,j}^k | z_{i,j})$$

assume $z_{i,j} \sim \text{Bernoulli}(\phi)$

$$\phi = \phi_i(z_{i,j})$$

parameters

$$\{P(x_{i,j}^k | z_{i,j})\}_{k=1}^2 = \{P(x_{i,j}^k | z_{i,j})\}_{k=1}^2$$

$$\frac{(c_m - 1) \sum_{i=1}^m (1 - w_{i,j})}{\sum_{i=1}^m (1 - w_{i,j})} : \phi = \phi_i(z_{i,j})$$

~~Singular: determinant is zero~~

$$\mathbf{Z} = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \mathbf{\mu})(\mathbf{x}_i - \mathbf{\mu})^T$$

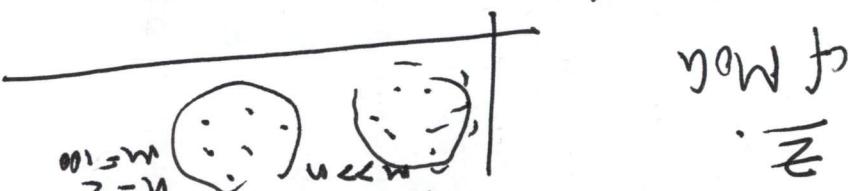
$$\mathbf{\mu} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i : \text{use MLE}$$

$$\mathbf{x} \sim N(\mathbf{\mu}, \mathbf{\Sigma}) \in \mathbb{R}^n$$

What can you do? Fit a Gaussian

$n >> m$: high dim \leftarrow 1000 dim, but only 20 measurements

$n \approx m$
What about:



But mathematically interesting: \therefore definition

Not really useful as formula, as NB says.

After simplifying model

How much do I think it is in the cluster.

M-step: no hard assignment, just

whether old class came from cluster or not

E-step: compute weights for each

The initial here

observed cluster

[Ng13]

So not a good model
Even though axes aligned
But you've thrown away all the correlation,

$$\mathbf{G} = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \mathbf{\mu})(\mathbf{x}_i - \mathbf{\mu})^T$$
$$\mathbf{Z} = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix} \quad \text{Eigenvalues}$$

Centroid \mathbf{Z} to be diagonal $\mathbf{x} \sim N(\mathbf{\mu}, \mathbf{Z})$

What to do?

\therefore Not a good model.

non-invertible:

$$((\mathbf{x}_1 - \mathbf{\mu})^T (\mathbf{x}_1 - \mathbf{\mu}), \dots, (\mathbf{x}_n - \mathbf{\mu})^T (\mathbf{x}_n - \mathbf{\mu}))$$

Counting of dimensions

infinitely thin
Gaussian

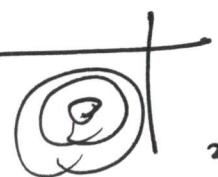
$n = m = 2$

That find the cov matrix will be singular, not full rank of non-invertible

(3)

④ In reality, you only see the x 's.
Informally, think of A as model
the data as coming from lower-dimensional
space plus noise.

FA affinities to model correct b/w variables

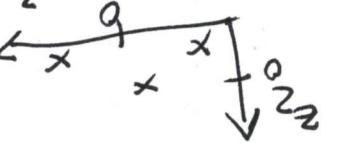


i.e. circular Gaussian

$$\text{O} \text{ covariance } S = \sigma^2 I = \begin{pmatrix} \sigma^2 & 0 & 0 \\ 0 & \sigma^2 & 0 \\ 0 & 0 & \sigma^2 \end{pmatrix}$$

Ng 13

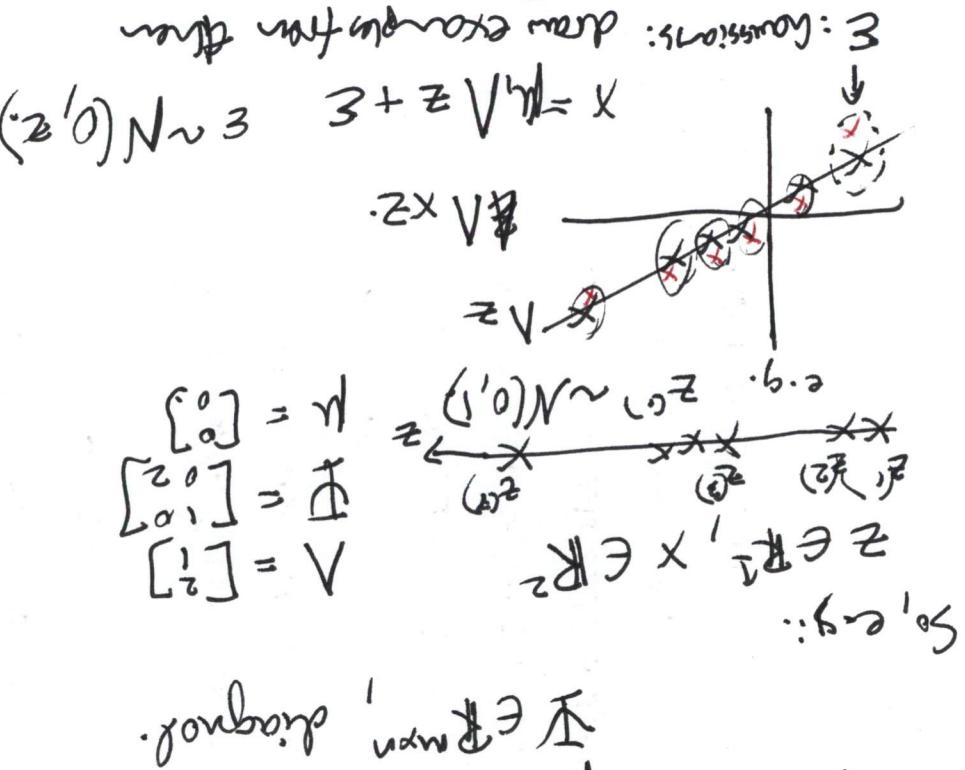
$\boxed{0:54:01}$
 $z \in \mathbb{R}^3, x \in \mathbb{R}^3$
So, example



the row! a dataset like a fat
matrix has more off the main diag
 $\mu + Vz$, where V is off the main diag

That's the model, how to fit the
factors.

PCA question?
- (short Gauss)
 $X = [x_1, x_2, \dots, x_n] \in \mathbb{R}^{n \times d}$
 $X = E[\mathbf{x}] + \mathbf{E}[e]$
partition vector



$$\begin{aligned} \Phi + I\lambda V &= \\ \left[1 \left(\lambda^2 - 2\lambda + 1 \right) (\lambda - 3 + 2V + V^T) \right] E &= 22E \\ 0 = I\lambda V &= \left[\begin{smallmatrix} 22 \\ 22 \end{smallmatrix} \right] \lambda V = \\ \left[\begin{smallmatrix} 22 \\ 22 \end{smallmatrix} \right] E - \left[\begin{smallmatrix} 22 \\ 22V \end{smallmatrix} \right] E &= \\ \left[\begin{smallmatrix} 22 \cdot (\lambda - 3 + 2V + V^T) \\ 22 - 2(\lambda - 3 + 2V + V^T) \end{smallmatrix} \right] E &= \\ E(\lambda - 3 + 2V + V^T) &= \\ E(\lambda - 3) &= \text{Co}(z) = I \\ \left[\begin{smallmatrix} (z - z)(x - x)^T \\ (z - z)(z - z)E \end{smallmatrix} \right] &= \\ \text{Co}(z) \text{ cov matrix: } E &= \end{aligned}$$

$(I'0)N \sim 3'3 + 2V + V^T = x'(I'0)N \sim 2$
 $(3'x^2N)N \sim \binom{x}{2}$

It turns out that you can go back.
 Under FA model

5

for λ_1, λ_2
 descent direction
 does not change
 $\dots = \lambda_1 \lambda_2$

$\lambda_1 \lambda_2 = \mu_1 + \mu_2$
 $(2\lambda_1 - 3) \sim 2\lambda_1 - 3 + \lambda_2$
 $(2\lambda_2 - 3) \sim 2\lambda_2 - 3 + \lambda_1$
 $\lambda_1 \lambda_2 \sim N(\mu_1, \sigma^2)$
 plug in the formula:

$$P(x_1 | x_2) = \frac{P(x_1, x_2)}{P(x_2)}$$

Conditional distribution

$x_1 \sim N(\mu_1, \sigma^2) \rightarrow$ marginal is a Gaussian

$$P(x_1) = \int P(x_1 | x_2) dx_2$$

$P(x)$

What are the marginal & conditional distributions?

Sufficiently a joint distribution

$\begin{bmatrix} E_{22} & E_{21} \\ E_{12} & E_{11} \end{bmatrix}$

$\begin{bmatrix} 22 \\ 22 \\ 213 \\ 113 \end{bmatrix} = \mathbf{B}$

$\leftrightarrow \leftrightarrow \leftrightarrow$

$\begin{bmatrix} 22 \\ 21 \\ 11 \end{bmatrix} = \mathbf{A}$

$\begin{bmatrix} 22 \\ 21 \\ 11 \end{bmatrix} = \mathbf{A}$
 $(3'x^2N)N \sim x$

Eq (3)

$$\text{Ex 3: } (x^2 + 1)^{-3} = (u^2)^{-3} \quad : \quad u = x$$

$$\text{Ansatz: } (\theta^2, x) \cdot P(z) = ((r^2)^2, \varphi) : \mathbb{A}^2 - \{(0,0)\}$$

∴ Use EM algorithm:

To take derivatives of θ with respect to α we can use the chain rule for MLE in this form.

$$(x)P(\omega)Z = (g)x \cdot M$$

$$(\pm \gamma_1 VV', \eta) N \sim x \quad P(x_{ij})$$

process to find passengers via

$$\left(\begin{bmatrix} I & N & V \\ 0 & I & V \\ 0 & 0 & I \end{bmatrix} \begin{bmatrix} w \\ y \\ 0 \end{bmatrix} \right) N \sim \begin{bmatrix} x \\ z \end{bmatrix} \quad \text{es } \boxed{Ng(\mathcal{B})N}$$

∴ $z_{(i)}$ is a continuous value var.
So it turns out that to input these E step.

$$\frac{(\theta^*)^{x_{(i)}}}{\theta^*(1-\theta^*)^{n-x_{(i)}}} \cdot \log(\theta^*) - \int_0^{\theta^*} \theta^* \ln \theta^* d\theta = \theta^*$$

$$f_{\text{EM}}(\theta^* | x_{(i)}) = P(z_{(i)} | x_{(i)}, \theta)$$

EM algo:
E step:

∴ Use EM algo

$$\theta^* = \frac{\partial \ell}{\partial \theta} : \text{MLE}$$

But calc + solve for MLE: θ^* is unsolvable

Likelihood of θ (posterior)

$$\prod_{i=1}^n P(x_{(i)} | \theta) = \prod_{i=1}^n \exp(-\frac{1}{2}(x_i - \mu)^2 / \sigma^2)$$

Given training set of parallel ex.

x used my model

$$= \text{marginal dist.} : (\frac{1}{\sqrt{2\pi}} \cdot \mathcal{N}(\mu, \sigma^2))$$

$P(x) = \prod_{i=1}^n P(x_{(i)})$

①

$$\left[z + \frac{1}{\sqrt{2\pi}} \cdot \left(\frac{x - \mu}{\sigma} \right) \right] \cdot \mathcal{N} \sim \mathcal{N}(\mu, \sigma^2)$$

↓ Gaussian draws and z using z

$$z = 1, n=2$$

parameters: $\mu \in \mathbb{R}^n$, $A \in \mathbb{R}^{n \times n}$ diagonal

$$z \sim \mathcal{N}(0, I)$$

$$z = \mu + \sqrt{z} + \epsilon : \text{constant} \times \text{gaussian}$$

$$z \sim \mathcal{N}(0, I) : \text{left of variable } z$$

$$\{x_{(i)}\}_{i=1}^n \in \mathbb{R}^n$$

EM: (max)

Principal component analysis (PCA)

factots analysis

29/6/22

Ng 14

$\frac{1}{2\pi} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \left(\log P(x_i) - \frac{1}{2} (x_i - \mu_i)^2 \right) dz_i$
 Now consider density
 $\sim N(\mu_i, \sigma_i^2)$
 Now
 $\frac{1}{2\pi} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \left(\log P(x_i) - \frac{1}{2} (x_i - \mu_i)^2 \right) dz_i$
 If you add up all terms
 $\sum_i \log P(x_i)$
 Then
 $\sum_i \log P(x_i) = \log \prod_i P(x_i)$

$$\sum_i \log P(x_i) = \sum_i \log \int_{-\infty}^{\infty} P(x_i | z_i) dz_i$$

For FA, only find the differences in parameters. So just maximize that one in
 $P(x_i | z_i)$.
 M-step: update

$$\begin{aligned}
 & \boxed{\sum_i \log P(z_i)} + \sum_i \log P(x_i | z_i) \\
 & \rightarrow \left[\theta_i | z_i \right] \log \left[\sum_j P(x_i | z_i) \right] \\
 & \left[\theta_i | z_i \right] \log \left[\sum_j P(x_i | z_i) \right] = \text{M-step}
 \end{aligned}$$

$E_{z_i \sim P_i} [\log P(x_i | z_i)] = \mu_i^T x_i$ result
 for simple first time is an expected value (a)

$$E_{z_i \sim P_i} [\log P(x_i | z_i)] = \mu_i^T x_i$$

$$(z_i - \mu_i)^T \Sigma_i^{-1} (z_i - \mu_i)$$

$$\Sigma_i^{-1} \mu_i$$

$$(a) write out \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-\frac{1}{2} (z_i - \mu_i)^T \Sigma_i^{-1} (z_i - \mu_i)}$$

$$\int_{-\infty}^{\infty} e^{-\frac{1}{2} (z_i - \mu_i)^T \Sigma_i^{-1} (z_i - \mu_i)} dz_i$$

key trick for computing this integral:

M-step: update

- key trick of mean & covariance

$\theta_i(z_i)$: compute vector of $\mu_i | z_i$

$$V_{i,-} (\bar{x} + W_{-i} V) - I = \sum_i \mu_i | z_i$$

$$(V_{i,-} | z_i) (z_i + V V) V - I = V_{i,-} | z_i$$

$$(V_{i,-} | z_i) \sim N(\mu_i | z_i, \Sigma_i^{-1})$$

$$\text{Principle of compensation and unbiasedness}$$

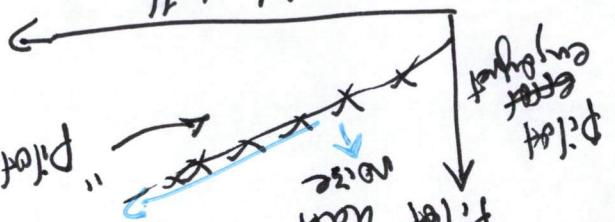
(3) $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$

Proof:

1. $E[\hat{\mu}] = E\left[\frac{1}{n} \sum_{i=1}^n x_i\right] = \frac{1}{n} \sum_{i=1}^n E[x_i] = \frac{1}{n} \sum_{i=1}^n \mu = \mu$
2. $\text{Var}(\hat{\mu}) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n x_i\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(x_i) = \frac{1}{n^2} \sum_{i=1}^n \sigma^2 = \frac{\sigma^2}{n}$
3. $\text{std } \hat{\mu} = \sqrt{\frac{\sigma^2}{n}}$

Here's the algorithm

Point SLL



e.g. 2. Now as we have seen in the helicoptrics
for now "Point of first" / "Point of second" / "Point of third" / ... / "Point of n-th".

e.g. 2. Now as we have seen in the helicoptrics
for now "Point of first" / "Point of second" / "Point of third" / ... / "Point of n-th".

Now if we want to find out what is the result of noise: Reduce the number of data points to get rid of noise.

Reduction of data points

What are they? x_1, x_2, \dots, x_n
Reducing it to k , then data set $(k \ll n)$
What are they? x_1, x_2, \dots, x_k

Effects are with x_i ~ Q:

$$V = \left(\sum_{i=1}^n (x_i - \mu) E[x_i] \right) \left(\sum_{i=1}^n E[(x_i - \mu)^2] \right)^{-1}$$

Ng 14

This is probably the most useful method I'll use.

I just wrote down a few steps of the derivation:

A: 33:40

Given that we can calculate the expected values of the random variable

$$E[x_i] = \mu_{x_i}$$

Substitute these into closure:

$$E[zz^T] = z(Ez)^T$$

$$E[zz^T] = E[(\mu_z + \epsilon_z)(\mu_z + \epsilon_z)^T]$$

$$E[zz^T] = \mu_z \mu_z^T + E[\epsilon_z \epsilon_z^T]$$

What are they?

$$E[\epsilon_z \epsilon_z^T] = \text{Cov}(\epsilon_z, \epsilon_z)$$

$$E[\epsilon_z \epsilon_z^T] = \text{Cov}(\epsilon_z, \epsilon_z)$$

Choose u_1, \dots, u_k to be k first eigenvectors of A

for $k < n$

to the best fit subspace

cover matrix Σ find principal eigenvector. This
To summarize: to find principal axis, construct

$$\text{① } \Delta u = \Sigma u - u \Sigma = 0 \quad \Delta \text{ construct subspace:}$$

orthogonal subspaces

$$(1 - u_1 u_1^T) u = u_1 \Sigma u - \Sigma u_1 u_1^T \leftarrow$$

$$= u_1 u_1^T \Sigma u + (1 - \|u\|)^2 \Sigma u = 0 \quad \text{eigenvalue of } A$$

$$A u = \lambda u \quad \text{eigenvector of } A$$

$$\Sigma = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$$

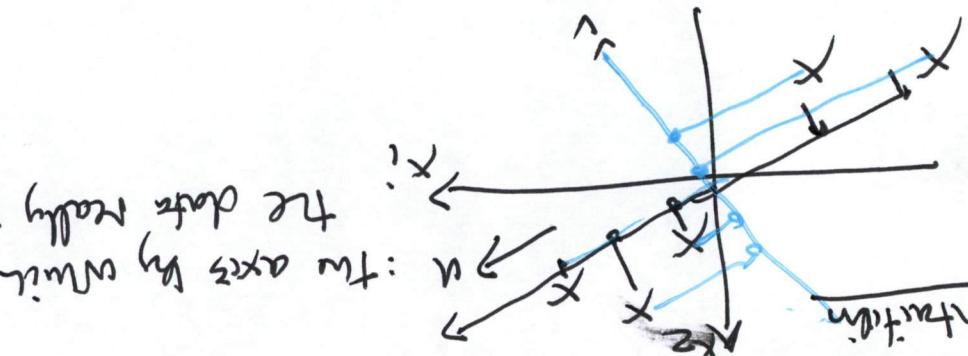
u must be Principal eigenvector of Σ

$$\|u\| = 1 \quad \text{and} \quad u = \frac{1}{\sqrt{n}} \sum_{i=1}^n x_i x_i^T u =$$

④

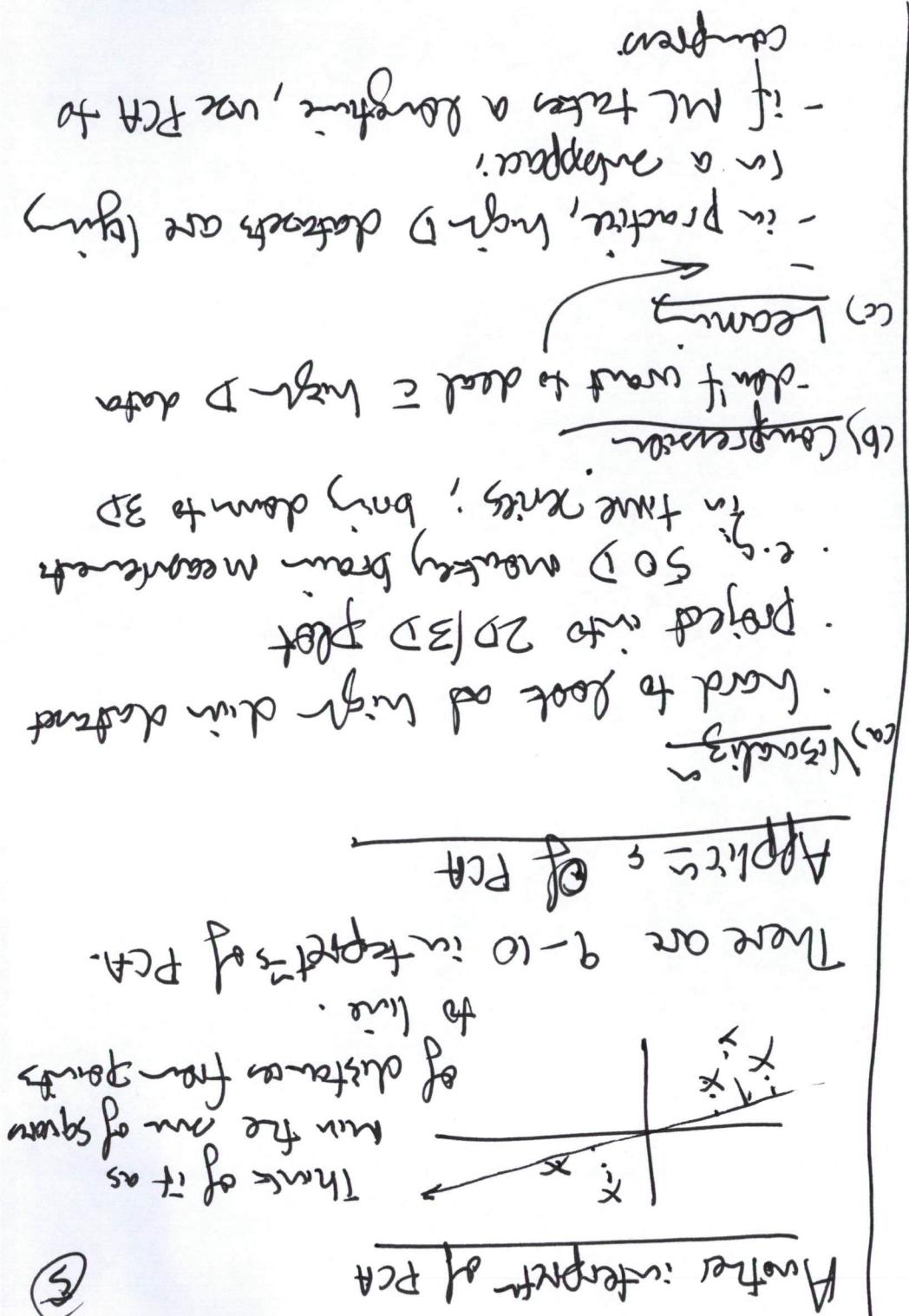
choose u : ~~$\frac{1}{\sqrt{n}} \sum_{i=1}^n (x_i x_i^T u)^2$~~
if $\|u\| = 1$, ~~$x_i^T u = 0$~~ perpendicular and u
Variance of Σ : find axis for the variance
(find the direction that allows greatest
variance: the direction is u)

The proxy of training set with u has



To find Principal axis of variance:

Principle axes = different scales



What would happen if different data come closely together? Loss of the benefit of SVD of course. tiny numerical issues for repeated data; sometimes dangerous to take freely in subspace. Think of u_i : some features there is ambiguity in u_i . It can be minus definite:

If some features are defective, but z_i is significant eigenvalues;

What do you mean by "principal eigenvalues"?

$u_i = \{u_{i1}, \dots, u_{ik}\}$ basis

$y^{(i)} = (u_1^T x^{(i)}, u_2^T x^{(i)}, \dots, u_k^T x^{(i)})$ $y^{(i)} \in \mathbb{R}^k$

Here $x^{(i)} \in \mathbb{R}^n$, new representation of data

∴ A small u_3

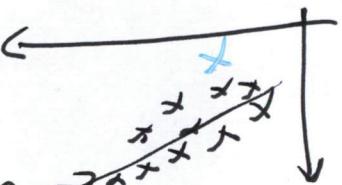
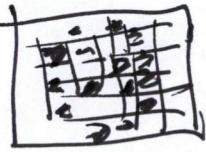
most data in the plane of the basis, but some outliers are.

Ng 14

High D, more features, \rightarrow complex hypothesis
 Class. trying to separate
 classes. In industry, this is done by hand.
 Feature selection
 to use select, auto-select
 like (sandy, flatland, etc.)
 lots of diff. e.g. weather
 : e.g. features
 $x_{(i)} \in R^{10^3}$
 linear
 $x_{(i)} \in R^{10^3}$
 etc.

PCA is useful to know. I use it and
 more often than it should be. Try to use
 the training data without changing
 the test set.
 SVD instead
 PCA related, but different
 - end for unstructured learning.

PCA application
 (a) Outliers (anomaly detection)
 - find suspicious for data; the root for
 finding fraud data + fit.
 - not good for this
 (b) Dimensionality reduction
 e.g. face recognition 100×100
 $x_{(i)} \in R^{10^3}$ shows appearance
 most of data lies on a
 subspace of low dim
 Use PCA to detect SVD subspace first
 give you a face! find other face that
 has same features.



Ng

Perception of PCA

- Difficulties
 - No guarantee of zero, unit variance
 - Find top k eigenvectors of $\Sigma \cdot Z = \Sigma = \sum_{i=1}^m x_i x_i^T$
 - Find large singular value

Application to test data: "LSI"

The difficulty of doing PCA, is that consider becomes large: $S = \text{Tr } \Sigma$

Test on high-dimensional vector

$\begin{bmatrix} \vdots & \vdots \\ 0 & a_{\text{middle}} \\ \vdots & \vdots \\ 1 & 1 \end{bmatrix}$ is zip

In LSI, we usually skip pre-processing step: don't scale up the weights of normalized, it would scale up the weights of raw words.

Problem: how similar are two documents $x^{(1)}, x^{(2)}$

One measure is $\sum_{i=1}^n (x_{i1}, x_{i2}) \cdot \cos \theta$. If θ is small, they are similar

$\sum_{i=1}^n (x_{i1}, x_{i2}) \stackrel{?}{=} \frac{\sum_{i=1}^n x_{i1} x_{i2}}{\|x^{(1)}\| \|x^{(2)}\|} = \cos \theta$

- Latent semantic analysis (LSA)
- Singular Value Decomposition
- Independent components analysis

- Latent semantic analysis (LSA)

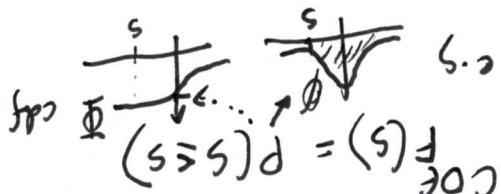
25

PCA

51 N

$s \hookrightarrow$ specifying $\beta_3(s)$ or specifying $f(s)$

$$f(s) = \int_{-\infty}^s p(t) dt$$



Random variable s , density $P(s)$

PDF - cumulative dist f \leq

Find independent components not principal
Shows slides of cortical party people, if audio
features of mixed sound

ICA

Avoid using PCA to reduce dim!

	K-means	Clustering / groups	Naive Bayes	PCA	Spectral	Feature	Model density not probability	$P(x)$	Model density not probability
--	---------	---------------------	-------------	-----	----------	---------	-------------------------------	--------	-------------------------------

CONTEXT

Naive Bayes: use SGD to compute eigenvectors in PCA

Nag 15

that way we know they cut things off.

- cutting off the 0 values, by convention

~~True~~ X is wide $[x_{\text{min}} \dots x_{\text{max}}] = [u \dots v]$

For 100 dimensions = can be computed directly

$X \in \mathbb{R}^{100 \times 50,000}$

top k columns of V acc top k eigenvectors of $X^T X = Z$

To get top k eigenvectors of Z :

$$\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} =$$

$Z = X^T X = X^T X = Z$: covariance matrix

$$\begin{bmatrix} -x_{(1)} & -x_{(2)} & -x_{(3)} \\ -x_{(2)} & -x_{(1)} & -x_{(3)} \\ -x_{(3)} & -x_{(2)} & -x_{(1)} \end{bmatrix} =$$

$$Z = \sum_i x_{(i)} x_{(i)}^T = \text{cov}(x)$$

"Design matrix"

$$\sum_{i=1}^n x_i = \sum_{i=1}^n y_i$$



$$([1,0], \eta \sim s) \quad \{1550\} \rightarrow (5, 2, 2)$$

$$P(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

$$d\omega_{\text{sys}} \propto s = P(s) \quad \text{and} \quad s = w - A_s = x$$

Dcurve of amalgamation: joint venture for expansion

-misses - Mrs

(I'0) $N \sim s^{\frac{1}{2}} S^{\frac{1}{2}}$ molar

missouri - now - (i) S.

permittimus, sis

—

Mathematics: Compt

152

for which use ≤ 6.3

$$\begin{bmatrix} \frac{1}{2}M & - \\ \frac{1}{2}m & - \\ -\frac{1}{2}m & - \\ -\frac{1}{2}M & - \end{bmatrix} = M \quad (1) \times M = (1) S$$

$$(1) \times M = (1)^M$$

load function

What are 's parts for

In microevolution each species adapts to its environment

$$x_{(i)} \in \mathbb{R}^n$$

The older & simpler x:

Speaker

A hand-drawn diagram of a wave function. It consists of a horizontal line with several vertical oscillations. The troughs of these oscillations are labeled with the letter 'S' in parentheses, indicating nodes where the wave amplitude is zero. The peaks are labeled with the letter 'A' in parentheses, indicating antinodes where the wave amplitude is at its maximum.

15

A handwritten note consisting of the letter 'S' followed by a horizontal line with a small arrow pointing to the left.

~~Perforated surfaces (speculars) (5)~~

just a safety check:

problems

(All have to be like it, e.g.)

$$\langle s \rangle, f = \langle s \rangle^f$$

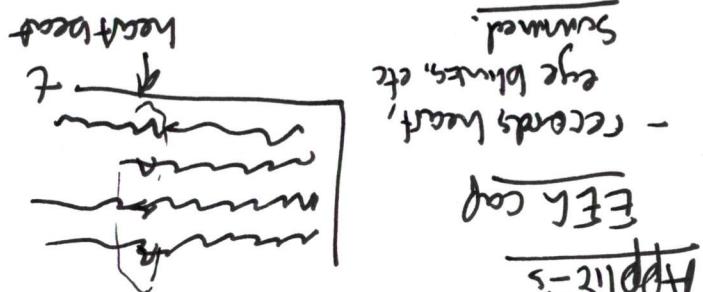
Batch processing of number of words of messages & words

Method of how brain processes information: CA

Image patterns make up neural message: CA

Hyper Hyperplane

- clean data
- remove affected coordinates & outliers
- eye blur
- find IC of brain signals - heart beats



$$S_{(i)} = W_{(i)} \rightarrow \text{separated out data series}$$

④

$$(W_T) \cdot \begin{bmatrix} (x_{(m)}) \\ \vdots \\ (x_{(1)}) \end{bmatrix} = (w) \Delta$$

$$\text{use sigmoid grad func: } \phi(s) = E(s)$$

$$|W| \cdot \left(\prod_{i=1}^n P_s(w_i x_i) \right) = \prod_{i=1}^n \underbrace{x_{(i)}}_{\text{data}} \cdots x_{(m)}$$

also works fine

> Also use $P(s) = e^{-s}$: Laplace dist will

not Gaussian

for $Q = I, \beta = 13$

$$\text{sigmoid: } F(s) = \frac{1+e^{-s}}{1-e^{-s}}$$

- can't be Gaussian Cdf, so convolution of m

$$\text{Conv } P_s(s) = ?$$

$$\left[\dots \begin{smallmatrix} \dots & \dots \\ \dots & \dots \end{smallmatrix} \dots \right] = M = H^{-1}$$

$$P(x) = \prod_{i=1}^n P_s(w_i x_i) \cdot |W|$$

$$P(s) = \prod_{i=1}^n P_s(s_i) \quad (s_i \text{ is input})$$

CA model

15/15

Ng 16

Reinforcement Learning [23, S6]

- MDPs: Markov Decision Process
- Value functions
- Value iteration
- Policy iteration

Reinforcement Learning
[23, S6]

Context
For supervised learning, we had the right answer
for unstructured learning, we had the right action
In between, some supervision: reinforcement
Picture of autonomous helicopter from Lecture 1:
- the football fields of stadium, heli: cargo
- difficult aerobatic maneuver.
Is difficult upskill down
Learning problem: ~~initial~~ - poss. of helicopter
- difficult aerobatic maneuver.

We don't know what the legal action is -
no training set for it. Instead, we reward
a good or bad action how well it is doing

$$\begin{aligned} P &: \text{reward function } R: S \rightarrow \mathbb{R} \\ \delta &: \text{discount factor } 0 \leq \delta \leq 1 \\ c(s) &\leq \sum p_{s,a}(s') = 1, p_{s,a}(s') \geq 0 \\ p_{s,a} &: \text{state transition distribution} \\ A &: \text{set of actions} \\ S &: \text{set of states} \end{aligned}$$

MDP $(S, A, \{p_{s,a}\}, \delta, R)$ S-tuple

MDP - makes decisions freely

- objectives encapsulated as series of steps!
forward this note

which moves were right, & which were
corrected: problems: you'd do if knew

Harder from optimised decisions, because no
unique right answer. You have to keep moving
e.g. 2: cars as obstacle. - like training a dog

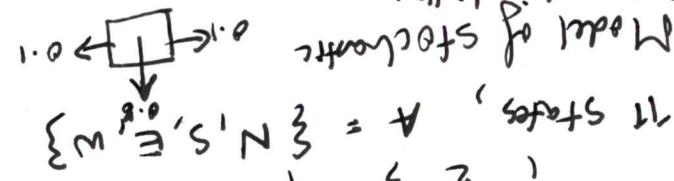
Example of an MDP

Ng 16

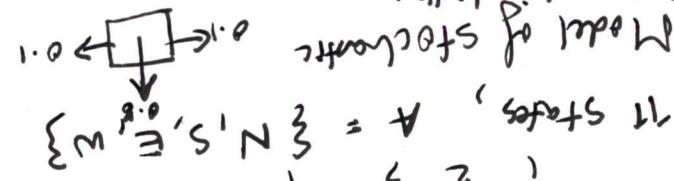
	1	2	3	4
1				
2	-1			
3	+1			

$$(4, 3)$$

- repeat frequency to recall top
adapted from P. Norvig



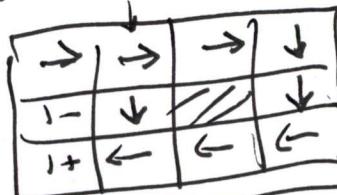
$$1 \text{ states}, A = \{N, S, E, W\}$$



$$\begin{aligned} P(3,1) N(3,3) &= 0 \\ P(3,1) N(2,1) &= 0.1 \\ P(3,0) N(4,1) &= 0.1 \\ P(3,1) N(3,2) &= 0.8 \\ P(3,0) N(2,1) &= 0.1 \\ P(3,0) N(4,1) &= 0.1 \\ P(3,1) N(3,3) &= 0 \\ P(4,1) &= -1 \\ R((4,3)) &= + \\ \text{Reward function:} \\ R((4,2)) &= -1 \\ R((4,1)) &= 0 \\ P(5,1) &= -0.02 \text{ for all other states} \\ R((5,2)) &= \end{aligned}$$

"The objective is the end of the problem: a 12th state!"
"A goal-conditioned state."

Optimal policy (strategy)



Why soft here? To avoid very large

-1. short but risky;

$$\text{Policy } \pi : S \leftarrow A : \text{for every state,}$$

$$E[\sum_{t=0}^{\infty} \gamma^t R(S_t) + \gamma R(S_t) + \gamma^2 R(S_2) + \dots]$$

total payoff

Goal is to choose actions over time that make
new runs.

So if diminishes the reward of each timestep
-0.5 < 1. Wins in the future weight less than
losses.

$R(S_0) + \gamma R(S_1) + \gamma^2 R(S_2) + \dots$: total payoff
After a while, we want/can apply R to $S_0 \dots$

Get to $S_2 \sim P_{S_0, a_0}$

Choose a_0

Get to $S_1 \sim P_{S_0, a_0}$

Get to $S_0 \sim P_{S_0, a_0}$

At state S_0

MDP walks like this:

②

Let's say this is our state space:
 Bellman's eqn gives a way of solving for V_{π}
 in closed form. If impeded linear constraints
 the value for π . The linear option of eqn's can
 be solved for the value for.
 $N = \{(3,1)\}$
 $V_{\pi} = 0.8V_{\pi}(3,1) + 0.1V_{\pi}(3,2) + 0.1V_{\pi}(4,1)$
 e.g. $V_{\pi}(3,1) = R(3,1) + 0.1V_{\pi}(1,0) + 0.1V_{\pi}(2,1)$
 write eqn for each state in MDP
 since are all values, || eqn!
 $\max_{\pi} V_{\pi}(s) = V_{\pi}(s) + \max_{\pi} \sum_{s'} P(s'|s, \pi) R(s', \pi) + \gamma V_{\pi}(s')$
 Bellman eqn for $V_{\pi}(s)$ rather than $V_{\pi}(s)$
 Optimal value function
 $V_{\pi}(s) = \max_{\pi} V_{\pi}(s)$
 more over all actions π
 expected payoff
 \downarrow

Used in Solving MDP
 Bellman's equation
 Previous (i.e. if action a from s is used
 action $\pi(s)$)
 $V_{\pi}(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi) V_{\pi}(s')$
 $\boxed{V_{\pi}(s) = E[R(s) + \gamma \sum_{s'} P(s'|s, \pi) V_{\pi}(s')]} \quad \Downarrow$

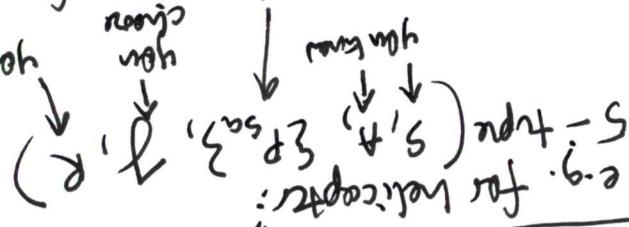
	58.	-1.0	
-1	28.		90.
1+	110.	23.	25.

	↓	↓	←	←
↓	↓	←	↑	
↓	↓	←	←	←

$V_{\pi}(s) = E[R(s) + \gamma \sum_{s'} P(s') V_{\pi}(s')] \quad \Downarrow$
 State s of execute π
 s.t. $V_{\pi}(s)$ is expected total payoff strategy π
 for every π , define value function V_{π} : $\pi \leftarrow S \mapsto V_{\pi}$
 optimal policy

Some definitions:
 Suboptimal
 If turns end that MDPs are good at solving
 ③

usually get more from others



What if you don't know P_{sa} ?

Value iteration

The computational efficiency step is solving the linear equations in Bellman's eqns. for millions of states, many equations. So we use Bellman-Bellman update:

$$(s) V(s) = \max_{a} \sum_{s'} P_{sa}(s) V(s')$$

$$- \text{left } V = V \quad (\text{so Bellman's eqn})$$

- repeat:

- until stable \Rightarrow (convergence)

Other standard algorithm: Policy iteration

e.g. expected total payoff is higher if you go left rather than right.

Using these numbers in to get V

14	15	16	17	18
1	69		28	
1	38	48	198	
1				

Run this algorithm on MDP

(a) synchronous update:
 $V_t(s) = V_{t-1}(s) + \gamma P(s'|s) V_t(s')$

2 usage of doing update:
 $V := V(N)$ - all states of one operator

This will make $V \leftarrow V(s)$: converge

$V(s) = P(s) + \max_{a} \sum_{s'} P(s'|s, a) V(s')$

Value iteration "iteration" \rightarrow array of elements

To reward the solution:

If I can compute V^* , then can get T^* , the compute - huge number of policies.

optimal policy. But $V^*(s)$ is not nice. Hard to

differentiate $\frac{\partial}{\partial V}$

$$(s) V(s) = \arg \max_a \sum_{s'} P_{sa}(s) V(s')$$

1:12:50

"Are you excited? Our first MDP algorithm"

- Update $V(s) = \text{argmax}_a P_{sa}(s) V(s)$ of V
- Some bellman's equations using Value iteration
- Update equations of P_{sa}
- Take action using π of q -value in MDP

Repet 3

Putting it together:

our way

more of a mix-and-match way of doing things!

$$\left(\frac{c}{\alpha} + \frac{1}{\beta} \right)$$

$$= \left(\frac{\# \text{times took action "a" in s}}{\# \text{times took action "a" in s}} \right)^{\frac{\alpha}{\beta}}$$

91