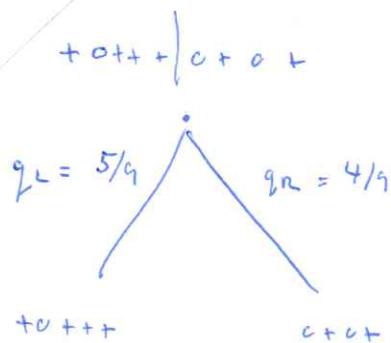


(62)

node of tree



$$q_L = (1/5, 4/5) \quad q_R = (1/2, 1/2)$$

Choose split that minimizes average impurity.

Say use H = entropy.

$$\text{Minimize } q_L H(q_L) + q_R H(q_R)$$

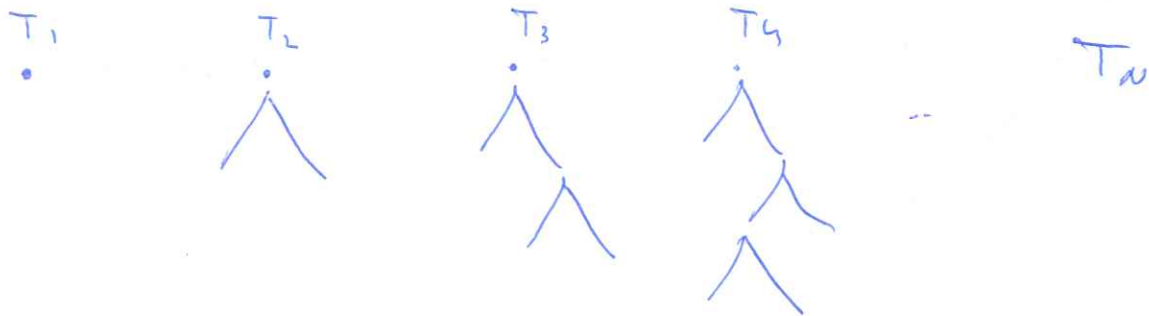
NB. Tree construction is greedy making locally optimal choices.

How deep to grow tree?

Tree-overfit.r

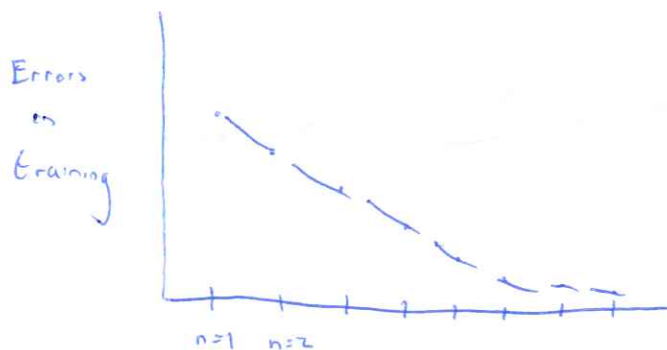
Overfitting

Let T_1, T_2, \dots, T_n be successive trees grown during training.

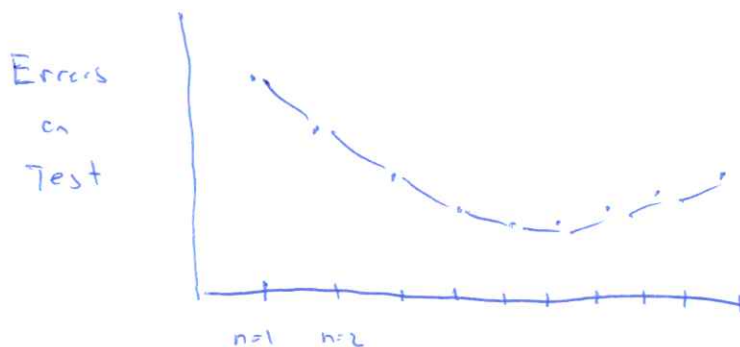


T_n has n terminal nodes ($|T_n| = n$)

The following graphs typical



More branches gives better result.



Performance improves for a while, then gets worse.

Overfitting

⊗ avoid overfitting?

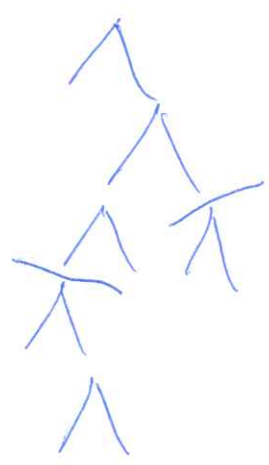
re complexity in choosing tree, T

Let $\alpha \in [0, \infty)$ be complexity param. Define

$$R_\alpha(T) = \underbrace{R(T)}_{\text{errors on training}} + \alpha \underbrace{|T|}_{\text{charge } \alpha \text{ for each split}}$$

Grow initial tree (deep) and consider rooted subtrees

$R_{\text{cut}} =$ all trees created by pruning initial tree.



Initial tree



$\in R_{\text{cut}}$

Define

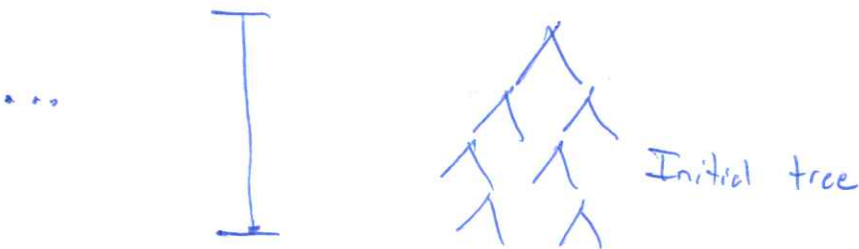
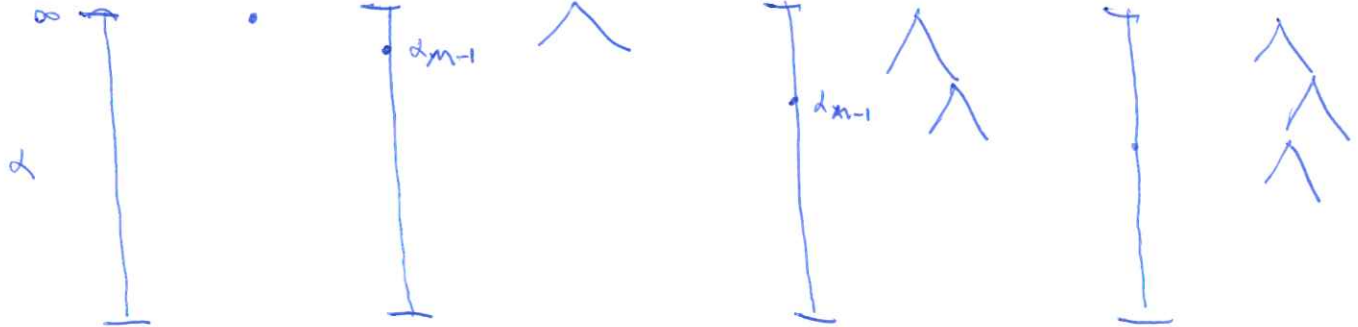
$$T_\alpha = \arg \min_{T \in R_{\text{cut}}} R_\alpha(T)$$

Note

① T_α is easy + cheap to construct

② The $\{T_\alpha\}$ are nested with $T_0 = \text{initial tree}$

$T_\infty = \text{null tree}$



Choose candidate complexity gains

$$\beta_1 = 0$$

$$\beta_2 = (\alpha_1 \alpha_2)^{1/2}$$

$$\beta_3 = (\alpha_2 \alpha_3)^{1/2}$$

\vdots

$$\beta_M = \infty$$

Select β by

Cross validation

Cross Validation

Have method for training classification or regression algorithm.

Divide training set G , randomly into G_1, \dots, G_N equal pieces.

$$\text{Let } \bar{G}_n = \bigcup_{j \neq n} G_j$$

CV estimates error rate of ~~the~~ training method.

$$\hat{\text{Error Rate}} = \frac{1}{N} \sum_{n=1}^N \text{Errors from } G_n \text{ when train on } \bar{G}_n$$

For CART let \hat{ER}_{β_m} be estimated ER using CV with complexity param β_m .

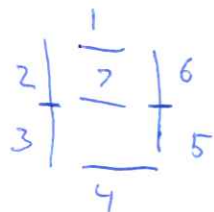
$$\beta^* = \arg \min_{\beta_m \in \{\beta_1, \dots, \beta_n\}} \hat{ER}_{\beta_m}$$

Now using all of G train tree with complexity param β^* .

~~the~~

Digit Recognition

Can get 10 digits $0, \dots, 9$ with 7 possible "strokes"



$$Z_0 = (1111110)$$

$$Z_1 = (0000110)$$

$$Z_2 = (1011011)$$

Create data as follows:

$$C_i \in \{0, \dots, 9\} \text{ at random}$$

$$X_{i1} \dots X_{i7} = Z_{C_i} \text{ with randomly flipped bits}$$

$$X_{i8} \dots X_{i24} = \text{noisy bits.}$$

Classify with digit-recognition.

"cosine" "trig" "brontescom"

$$\begin{pmatrix} \text{class 1} \\ \text{class 2} \\ \vdots \\ \text{class } k \end{pmatrix} \begin{pmatrix} q_{11} & \dots & q_{1D} \\ q_{21} & \dots & q_{2D} \\ \vdots & & \vdots \\ q_{k1} & \dots & q_{kD} \end{pmatrix} = \begin{matrix} Q \\ k \times D \end{matrix}$$

$$\begin{pmatrix} \text{obs 1} \\ \text{obs 2} \\ \vdots \\ \text{obs } n \end{pmatrix} \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1D} \\ x_{21} & x_{22} & \dots & x_{2D} \\ \vdots & & & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nD} \end{pmatrix} = \begin{matrix} X \\ n \times D \end{matrix}$$

$$a) \prod_{i=1}^n \sum_{k=1}^K \pi_k \prod_{d=1}^D q_{kd}^{x_{id}} (1 - q_{kd})^{1-x_{id}} = P(X | \pi, Q)$$


$$b) \delta_{ik} = P(\text{class} = k | x_i) = \frac{\pi_k \prod_{d=1}^D q_{kd}^{x_{id}} (1 - q_{kd})^{1-x_{id}}}{\sum_{k'} \uparrow}$$

$$c) \hat{\pi}_k = \frac{n_k}{n} \quad n_k = \sum_i \delta_{ik}$$

$$\hat{q}_{kd} = \frac{\text{count of class } k \text{ with } x_{id}=1}{\text{\# count of class } k} = \frac{\sum_{\{i: x_{id}=1\}} \delta_{ik}}{n_k} = \frac{\sum \delta_{ik} x_{id}}{n_k}$$

$$\begin{aligned}
 d) \quad \log q(z, X; \pi, Q) &= \log \prod_{i=1}^n \pi_{z_i} \prod_{d=1}^D q_{z,d}^{x_{i,d}} (1 - q_{z,d})^{1-x_{i,d}} \\
 &= \sum_{i=1}^n \log \pi_{z_i} + \sum_{d=1}^D x_{i,d} \log q_{z,d} + (1-x_{i,d}) \log (1 - q_{z,d})
 \end{aligned}$$

$$\begin{aligned}
 e) \quad \mathbb{E}_{\pi^{\text{old}}, Q^{\text{old}}} \log q(z, X; \pi, Q) &= \\
 \sum_{i=1}^n \sum_{k=1}^K \delta_{i,k} \log \pi_k &+ \sum_{d=1}^D x_{i,d} \sum_{k=1}^K \delta_{i,k} \log q_{k,d} + (1-x_{i,d}) \sum_{k=1}^K \delta_{i,k} \log (1 - q_{k,d})
 \end{aligned}$$



π_k only appear here

$$f) \quad \hat{\pi}^{\text{new}} = \arg \max_{\pi} \sum_{k=1}^K n_k \log \pi_k \quad \text{s.t. } \sum \pi_k = 1$$

$$\frac{n_k}{\pi_k} = \lambda \quad \Rightarrow \quad \frac{n_k = \lambda \pi_k}{n = \lambda}$$

$$\frac{n_k}{\pi_k} = n \quad \Rightarrow \quad \pi_k^{\text{new}} = \frac{n_k}{n} \quad (\text{as always})$$

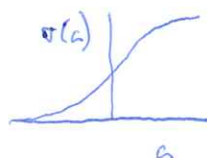
Neural Networks

Recall 2-class logistic regression. Have multivariate ~~vector~~ vector $x \in \mathbb{R}^D$

Model

$$P(\text{class} = 1 | x) = \sigma\left(\sum_{d=1}^D w_d x_d\right) = \sigma(w^T x)$$

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



Inputs

(x)

(x)

(x)

(x)

Output

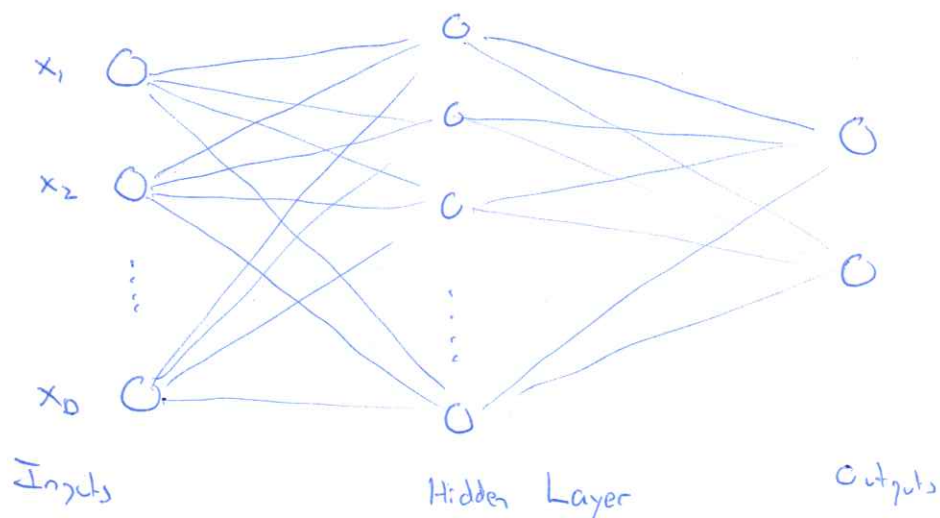
$$\sigma(a) = \sigma(w^T x)$$

$$a = \sum_{d=1}^D w_d x_d$$

Output is nonlinear xform of linear comb. of inputs.

Neural Network generalizes to multiple ^{logistic-type} nodes

Typical Network

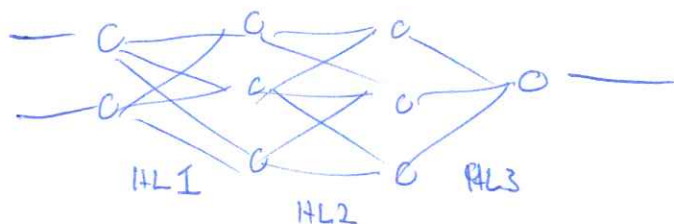


Let z_1, \dots, z_n be inputs to any node

Network Topology

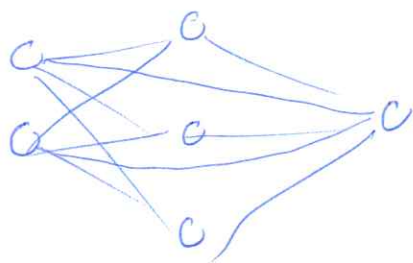
Can have

- 1) Any # input nodes, hidden nodes, outputs
- 2) Any # of hidden layers



- 3) Layers need not be completely connected, but must be "feed forward". That is, if $D(v_i)$ is the depth of vertex (node) v_i and $v_i \rightarrow v_j$ in network, then $D(v_j) > D(v_i)$

[Graph cannot have cycles]



Connections can skip over layers.

~~Non-linearity~~ Activation

At each hidden node v_j , the activation $a_j = \sum_{i=1}^D w_{ji} z_i$
 where $\{z_i\}$ are inputs. (First layer has inputs $\{x_i\}$)

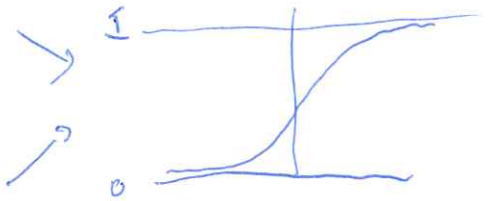
The activation is governed by "weights" w_{ji} . These are parameters of nn . (Learned during training)

Nonlinearity

At node v_j have output $z_j = h(a_j) = h\left(\sum_{i=0}^n w_{ji} z_i\right) = h(w_j^t z)$

Common choices for h :

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



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

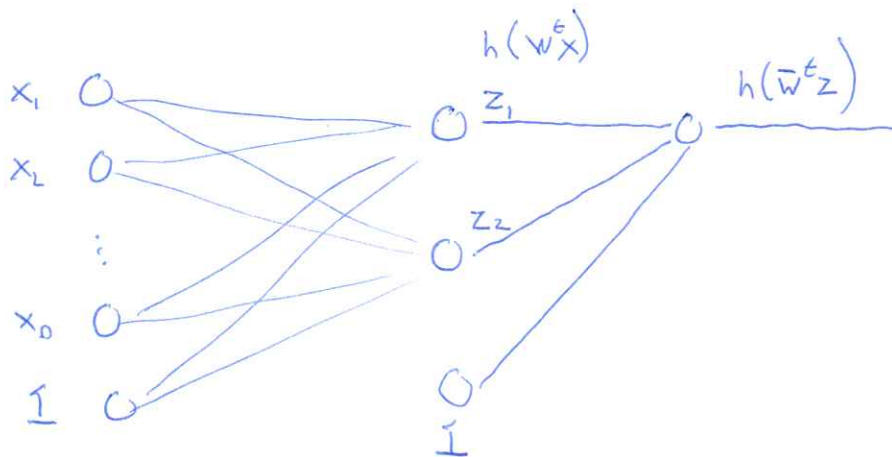


Bias

Sometimes activation written as

$$a_j = w_j^t z + \underbrace{w_{j0}}_{\text{bias term}}$$

As in regression can add 1 to input



Then outputs can be written $h(w^t x)$ or $h(w^t z)$

Output of Neural Network

Form of output depends on goal of network

3 cases: Regression (continuous output)

2-class classification

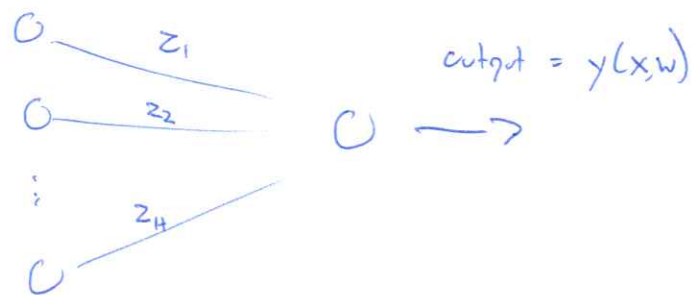
K-class classification

① Regression (continuous output [predict lifetime of object])

Denote output of NN by $y(x, w)$

\nwarrow inputs \nearrow weights

Let z_1, \dots, z_H be inputs to final node



Model
$$y(x, w) = \sum_{i=1}^H w_i z_i = w^t z$$

NB No non linear (h) trans on output.

② 2-class classification

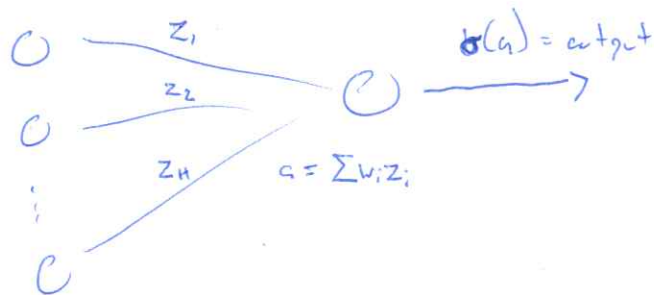
z_1, \dots, z_H inputs to final node

Output: $y(x, w) = P(C=1 | x, w)$

Model

$$y(x, w) = \sigma\left(\sum_i w_i z_i\right) = \sigma(w^t z)$$

$$\sigma(a) = \frac{1}{1 + e^{-a}} \quad \text{as in logistic regression.}$$



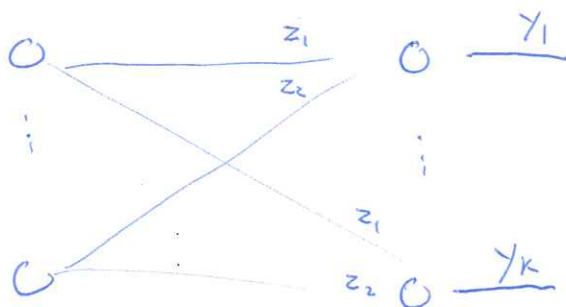
③ k-class classification

Output $y_k(x, w) = P(C=k | x, w)$

$$\sum_k y_k(x, w) = 1$$

$$0 \leq y_k(x, w) \leq 1$$

$$y_k(x, w) = \frac{e^{w_k^t z}}{\sum_{k'=1}^k e^{w_{k'}^t z}} \quad (\text{softmax})$$



Universal Approximation

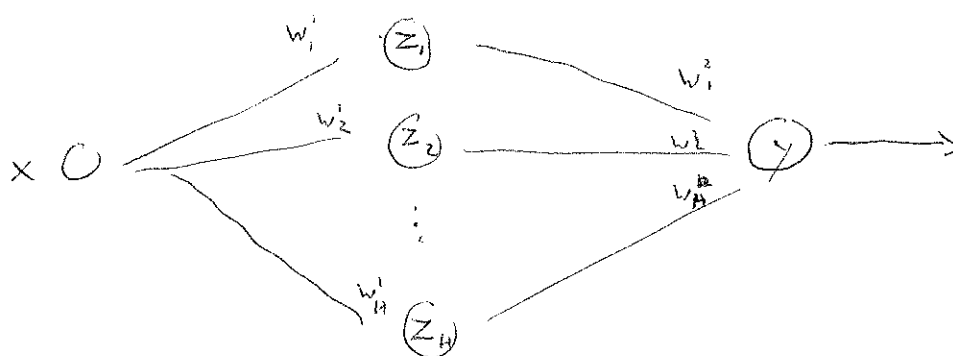
Neural net $y(x, w)$ represents function.

\swarrow input
 \searrow output

Loosely speaking, NN can approx any function with arbitrary precision.

More precisely, suppose $x \in [0, 1]$ with $t(x)$ "target" fn (what we want to approx.)

Let $y(x, w)$ be NN with H hidden units as follows.



$$y(x, w) = y(x, w_1^1, \dots, w_H^1, w_1^2, \dots, w_H^2) = \sum_{j=1}^H w_j^2 h(w_j^1 x)$$

Then

$$y(x, w) \xrightarrow{H \rightarrow \infty} t(x)$$

if weights $\{w_1^1, \dots, w_H^1, w_1^2, \dots, w_H^2\}$ chosen properly.

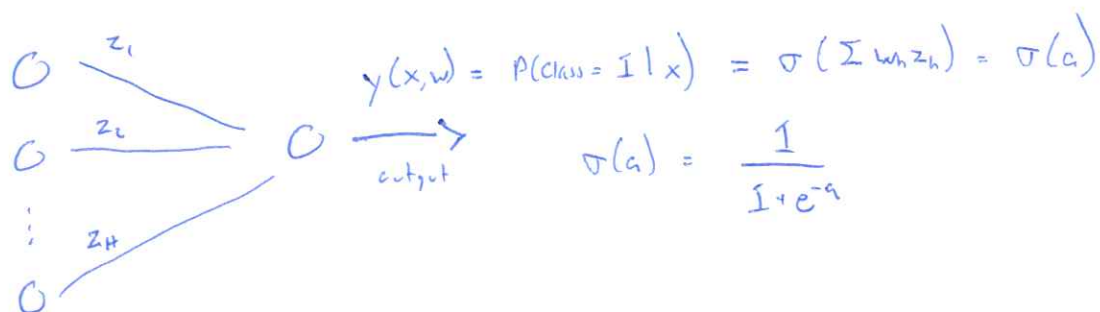
Note Other methods are also universal approximators

① GMM can approx any dist as # of mixture comp (k) increases

② CART used for regression (as an exam) can approx any fn.

Training the Network

Consider 2-class classification

Have data $(x_1, \epsilon_1), (x_2, \epsilon_2), \dots, (x_n, \epsilon_n)$ $\epsilon_i \in \{0, 1\}$

$$P(\epsilon_1, \dots, \epsilon_n | x_1, \dots, x_n, w) = \prod_{i=1}^n P(\text{class} = 1 | x_i, w)^{\epsilon_i} (1 - P(\text{class} = 1 | x_i, w))^{1-\epsilon_i}$$

$$= \prod_{i=1}^n y(x_i, w)^{\epsilon_i} (1 - y(x_i, w))^{1-\epsilon_i}$$

$$E(w) = -\log P(\epsilon_1, \dots, \epsilon_n | x_1, \dots, x_n, w) = -\sum_{i=1}^n \epsilon_i \log y(x_i, w) + (1-\epsilon_i) \log (1-y(x_i, w))$$

$$= \sum_{i=1}^n E_i(w)$$

$$\frac{\partial E_i(w)}{\partial a} = \frac{-\epsilon_i}{y(x_i, w)} \frac{\partial y(x_i, w)}{\partial a} + \frac{1-\epsilon_i}{1-y(x_i, w)} \frac{\partial y(x_i, w)}{\partial a}$$

$$= \left(\frac{1-\epsilon_i}{1-y(x_i, w)} - \frac{\epsilon_i}{y(x_i, w)} \right) y(x_i, w) (1-y(x_i, w))$$

$$= (1-\epsilon_i) y(x_i, w) - \epsilon_i (1-y(x_i, w))$$

$$= y(x_i, w) - \epsilon_i$$

$$= y_i - \epsilon_i$$

Gradient Descent by Backpropagation

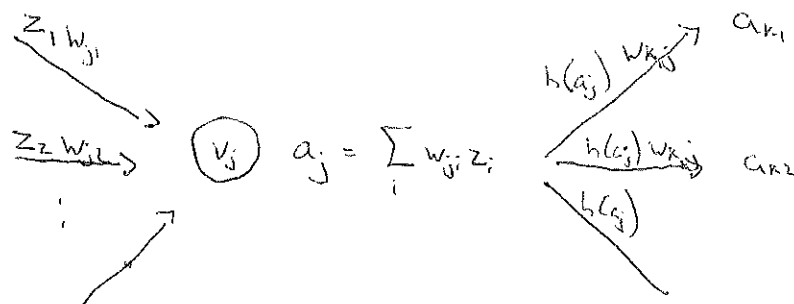
76

NN parametrized by weights: w $y(x, w)$

Want to minimize error $E(w)$

$$E(w) = \sum_{n=1}^N \underbrace{E_n(w)}_{\text{error for single observation}} \Rightarrow \nabla_w E(w) = \sum_{n=1}^N \nabla_w E_n(w)$$

Will focus on $\nabla E_n(w) = \left\{ \frac{\partial E_n}{\partial w_{ji}} \right\}_{\text{all connections } i \rightarrow j}$



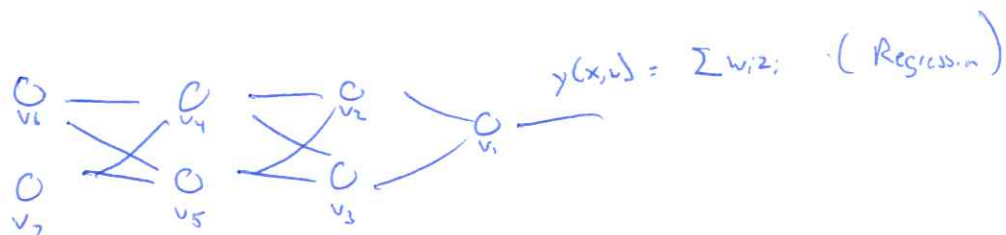
E_n depends on w_{ji} only through a_j

$$\left[\begin{aligned} E_n(w_{ji}, \dots) &= E_n(a_j(w_{ji}, \dots), \dots) \\ \frac{\partial E_n}{\partial w_{ji}} &= \underbrace{\left(\frac{\partial E_n}{\partial a_j} \right)}_{\delta_j} \frac{\partial a_j}{\partial w_{ji}} = \delta_j z_i \end{aligned} \right] \quad \begin{aligned} E_n(a_j, \dots) &= E_n(a_{k1}(a_j), a_{k2}(a_j), \dots) \\ E_n &\text{ depends on } a_j \text{ only through } \{a_k: j \rightarrow k\} \end{aligned}$$

$$\delta_j = \frac{\partial E_n}{\partial a_j} = \sum_{k: j \rightarrow k} \frac{\partial E_n}{\partial a_k} \frac{\partial a_k}{\partial a_j} = \sum_{k: j \rightarrow k} \delta_k w_{kj} h'(a_j) = h'(a_j) \sum_{k: j \rightarrow k} \delta_k w_{kj}$$

NB We develop δ_j in terms nodes "downstream" of v_j . This is backpropagation.

An Example



$$\delta_1 = \frac{\partial E_n}{\partial a_1} = (y_n - t_n)$$

$$\delta_2 = w_{12} \delta_1 h'(a_2)$$

$$\delta_3 = v_{13} \delta_1 h'(a_3)$$

$$\delta_4 = (v_{24} \delta_2 + w_{34} \delta_3) h'(a_4)$$

$$\delta_5 = (v_{25} \delta_2 + v_{35} \delta_3) h'(a_5)$$

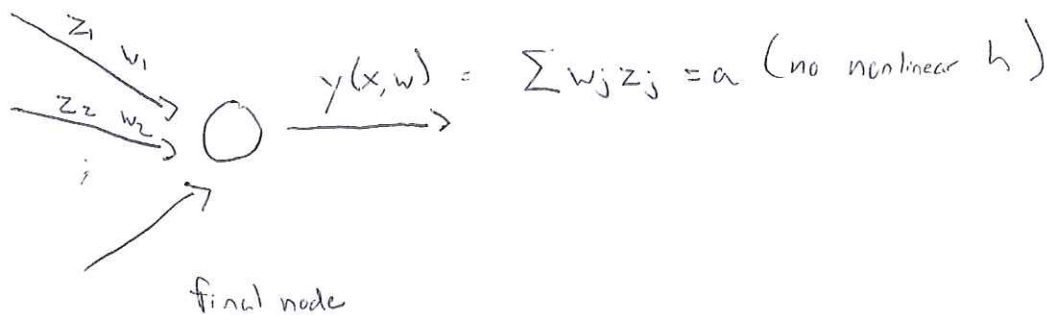
$$\frac{\partial E_n}{\partial w_{12}} = z_2 \delta_1$$

$$\frac{\partial E_n}{\partial v_{13}} = z_3 \delta_1$$

$$\frac{\partial E_n}{\partial w_{24}} = z_4 \delta_2$$

$$\frac{\partial E_n}{\partial v_{46}} = x_6 \delta_4$$

Regression Case



For regression use sq error as ~~loss~~ ^{objective} fn. Have $(x_1, t_1) \dots (x_n, t_n)$

$$E(w) = \frac{1}{2} \sum_{n=1}^N (y(x_n, w) - t_n)^2$$

$$E_n(w) = \frac{1}{2} (y_n - t_n)^2$$

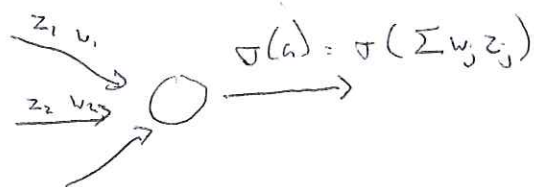
$$y_n = y(x_n, w)$$

$$E(w) = \frac{1}{2} \sum_{n=1}^N (y_n - t_n)^2 = \sum_{n=1}^N E_n(w)$$

As before differentiate E_n wrt. final activation a

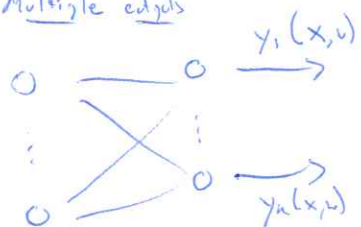
$$\frac{\partial E_n}{\partial a} = y_n - t_n$$

Note This is exactly what we get for 2-class classification



Even though had different objective fn. (log like) and different h fn.

Multiple edges



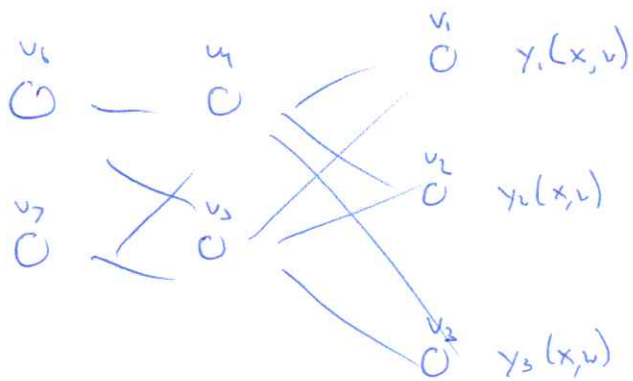
$$E(w) = \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K (y_k(x_n, w) - t_{nk})^2 = \sum_{n=1}^N E_n(w)$$

$$E_n(w) = \frac{1}{2} \sum_{k=1}^K (y_k(x_n, w) - t_{nk})^2$$

$$\frac{\partial E_n}{\partial a_k} = y_k(x_n, w) - t_{nk}$$

Regression E_T

75



$$\delta_1 = y_1(x_1, w) - t_{1n}$$

$$\frac{\partial E_T}{\partial w_{14}} = z_4 \delta_1$$

$$\delta_2 = y_2(x_1, v) - t_{2n}$$

$$\delta_3 = y_3(x_1, v) - t_{3n}$$

$$\frac{\partial E_T}{\partial w_{46}} = x_6 \delta_4$$

$$\delta_4 = h'(a_4) \sum_{k=1}^3 w_{k4} \delta_k$$

$$\delta_5 = h'(a_5) \sum_{k=1}^3 w_{k5} \delta_k$$

Regularization of Neural Nets

NN's can have many hidden nodes (= many parameters (= weights))

This can lead to overfitting.

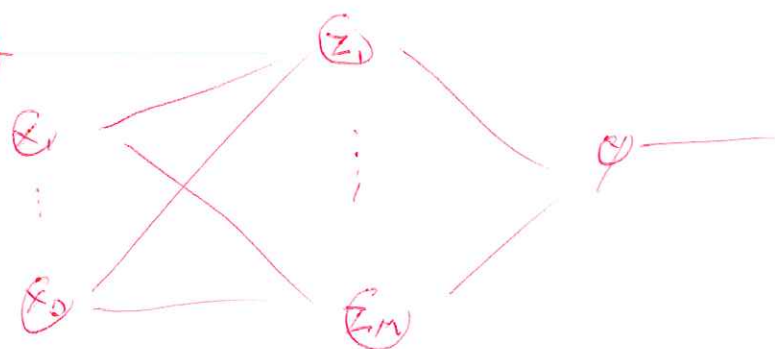
How to address this?

① Cross Validations

Have separate validation set $(x_i^v, t_i^v) \dots$ not used in training.

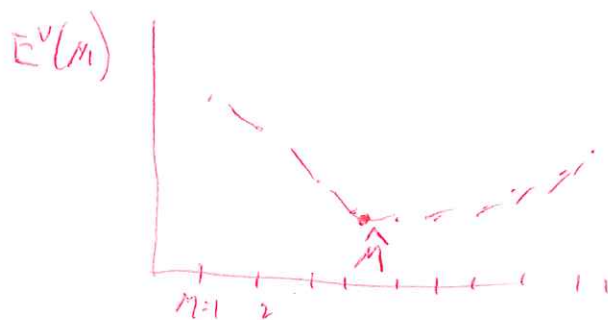
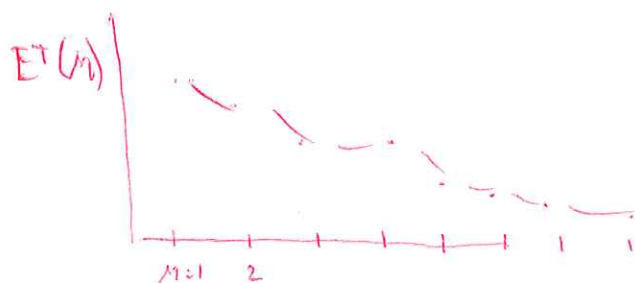
Suppose network has M hidden nodes with known topology (connectivity structure) but M unknown.

Eg



For $M = 1, 2, \dots$ train network using backprop.

Let $E^T(M)$ be error on training set after training with M .



Typically $E^T(M)$ decreasing in M .

Let $E^V(M)$ be error on validation after training.

Choose \hat{M}^* by $\hat{M} = \arg \min_M E^V(M)$

(2) Early Stopping

Since NN training is gradient descent, training iterative.

Let $y(x, w^n)$ $n=1, 2, \dots$

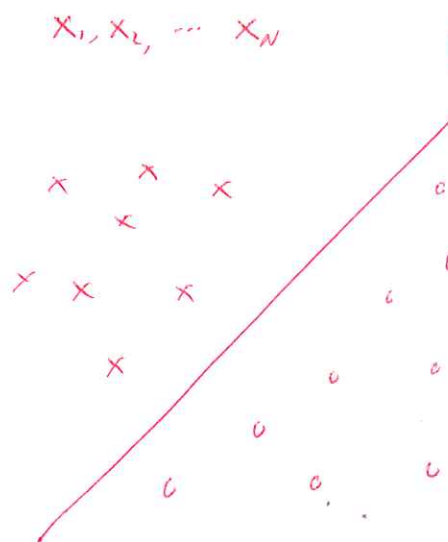
be network (function estimate) after n iterations of grad. descent.

Let $E^v(n)$ be error on separate validation set.

Stop training when $E^v(n)$ starts to increase.

Support Vector Machines

Consider 2-class classification problem on D -dimensional data



$$\{x: w^T x + b = 0\}$$

If $D=2$ (as in figure) and $w \in \mathbb{R}^D, b \in \mathbb{R}$

$\{x: w^T x + b = 0\}$ is a line

If $D=3$

$\{x: w^T x + b = 0\}$ is a plane

If $D > 3$

$\{x: w^T x + b = 0\}$ is a "hyperplane".

~~We seek $w \in \mathbb{R}^D, b \in \mathbb{R}$ s.t.~~

Can view w, b as classifier. Suppose classes are $\{+1, -1\}$

If $w^T x + b > 0$ classify as $+1$

If $w^T x + b < 0$ classify as -1

We seek w, b that perfectly classify x_1, x_2, \dots, x_N . That is
 Suppose classes are t_1, \dots, t_N $t_n \in \{-1, +1\}$. Seek w, b s.t.

$$t_n (w^T x_n + b) > 0 \quad n = 1 \dots N$$

Note:

May not be possible to classify perfectly, but can consider xform

$x \rightarrow \phi(x)$ (feature vector) s.t.

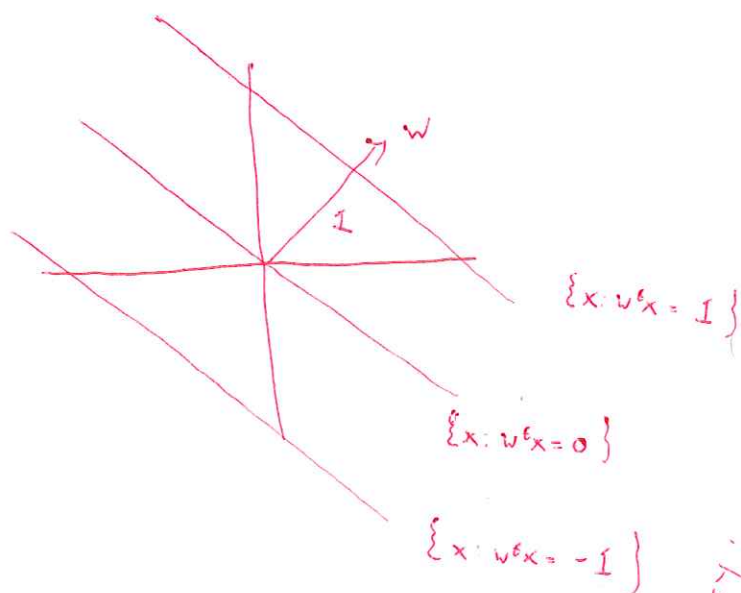
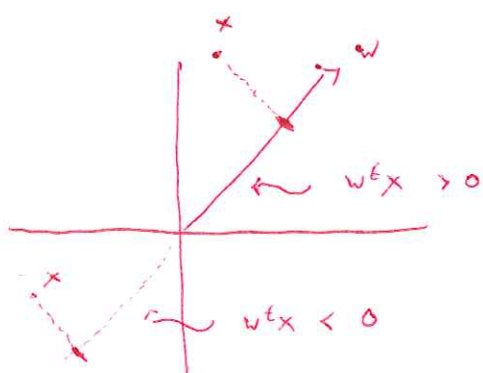
$$t_n (w^T \phi(x_n) + b) > 0 \quad n = 1 \dots N$$

Geometry

$w^t x = \sum_{i=1}^D w_i x_i$ is inner product between w and x

$$(\dots w^t \dots) \begin{pmatrix} \vdots \\ x \\ \vdots \end{pmatrix}$$

If $\|w\| = \left(\sum w_i^2\right)^{1/2} = 1$ $w^t x$ gives signed length of proj of x onto w .



① Suppose $\{x: w^t x = 0\}$ is decision boundary
Then $|w^t x|$ gives distance of x from decision boundary

② Suppose $\{x: w^t x + b = 0\}$ is decision bdy.

Then $|w^t x + b|$ gives dist of x from decision bdy.

③ Have $(x_1, t_1) \dots (x_n, t_n)$ $t_n \in \{-1, +1\}$.

Suppose have w, b that perfectly classify $x_1 \dots x_n$.

That is, $t_n = \begin{cases} +1 & \text{if } w^t x_n + b > 0 \\ -1 & \text{if } w^t x_n + b < 0 \end{cases}$

Then $\forall x$
For $x_1 \dots x_n$ $t_n(w^t x_n + b) > 0$
so $t_n(w^t x_n + b)$ gives dist. to bdy.
(note no abs value!)

What if $\|w\| \neq 1$?

Can still classify according to sign of $w^t x + b$

$$\text{class}(x) = \begin{cases} +1 & w^t x + b > 0 \\ -1 & w^t x + b < 0 \end{cases}$$

Note that if $\alpha > 0$, scaling $w^t x + b$ by α gives equivalent classifier

$$\text{class}(x) = \begin{cases} +1 & \alpha(w^t x + b) > 0 \\ -1 & \alpha(w^t x + b) < 0 \end{cases}$$

Suppose we choose α s.t. $\|\alpha w\| = 1$. That is $\alpha = \frac{1}{\|w\|}$.

$$\text{Then } \|\alpha w\| = \left\| \frac{w}{\|w\|} \right\| = \frac{1}{\|w\|} \|w\| = 1$$

$$\text{class}(x) = \begin{cases} +1 & \frac{w^t}{\|w\|} x + \frac{b}{\|w\|} > 0 \\ -1 & \frac{w^t}{\|w\|} x + \frac{b}{\|w\|} < 0 \end{cases}$$

Thus the distance of x to bdry given by

$$\frac{w^t}{\|w\|} x + \frac{b}{\|w\|}$$

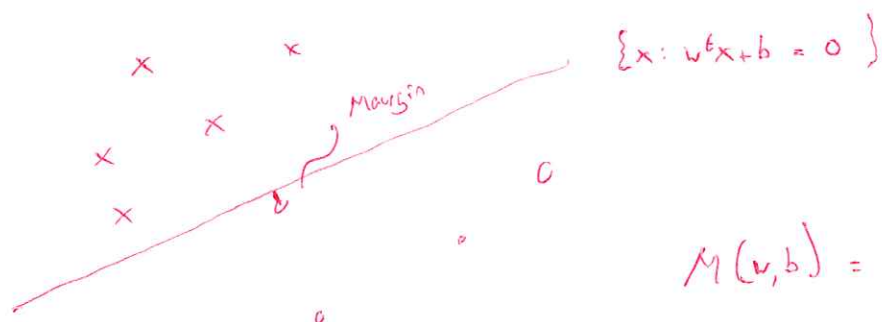
($\|w\| \neq 1$)

Thus if w, b correctly classify all x_1, \dots, x_n then dist. of x_n to bdry given by

$$L_n \left(\frac{w^t}{\|w\|} x_n + \frac{b}{\|w\|} \right) = L_n \left(\frac{w^t x_n + b}{\|w\|} \right)$$

Margin

If w, b classify all x_1, \dots, x_N correctly, the margin of w, b is the minimum dist to bdy



$$M(w, b) = \min_{n=1 \dots N} \epsilon_n \left(\frac{w^T x_n + b}{\|w\|} \right)$$

Support Vector Machine (SVM)

The SVM seeks w, b giving maximum margin.

$$(\hat{w}, \hat{b}) = \arg \max_{w, b} \min_n \epsilon_n \left(\frac{w^T x_n + b}{\|w\|} \right)$$

Consider xform $w \rightarrow \alpha w$
 $b \rightarrow \alpha b$ for some $\alpha > 1$.

Note that classifier is unchanged since $w^T x + b$ unchanged.

However, xform scales $w^T x + b$ by factor of α .

For any w, b that correctly classifies, choose α to give ~~margin~~ ϵ_n .

$$\min_n \epsilon_n(w^T x_n + b) = 1. \text{ Then}$$

$$(\hat{w}, \hat{b}) = \arg \max_{w, b} \min_n \frac{\epsilon_n(w^T x_n + b)}{\|w\|} \quad \text{s.t.} \quad \min_n \epsilon_n(w^T x_n + b) = 1$$

$$= \arg \max_{w, b} \frac{1}{\|w\|} \quad \text{s.t.} \quad \min_n \epsilon_n(w^T x_n + b) = 1$$

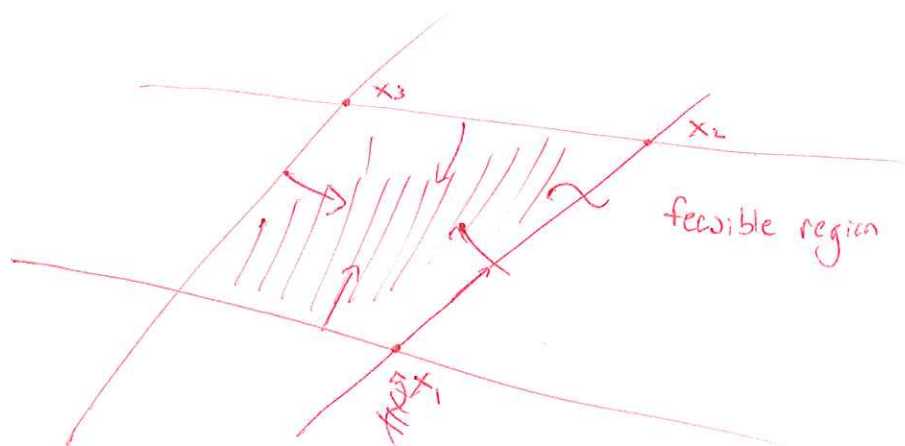
$$= \arg \max_{w, b} \frac{1}{\|w\|} \quad \text{s.t.} \quad \min_n \epsilon_n(w^T x_n + b) \geq 1$$

$$= \arg \max_{w, b} \frac{1}{\|w\|^2} = w^T w \quad \text{s.t.} \quad \min_n \epsilon_n(w^T x_n + b) \geq 1$$

This is a quadratic programming (QP) problem.

Quadratic Programming

Linear programming seeks to maximize linear fn subject to linear constraints.



Traverse feasible corners
until adjacent corners give
worse value of objective fn.

(Simplex method)

Quadratic Programming minimizes quadratic fn s.t. linear constraints.

Q.P.

Minimizes $x^T D x + d^T x$ s.t.

$$a_1^T x \geq c_1$$

$$a_2^T x \geq c_2$$

$$\vdots$$

$$a_n^T x \geq c_n$$

$$\begin{pmatrix} c_1 \\ \vdots \\ c_n \end{pmatrix} \uparrow \\ \Leftrightarrow A x \geq c$$

Our Problem

Minimize $\|w\|^2 = w^T w$ s.t. $\epsilon_n (w^T x_n + b) \geq 1 \quad n=1 \dots N$

$$\begin{pmatrix} -a_1^T & - \\ -a_2^T & - \\ \vdots & \vdots \\ -a_n^T & - \end{pmatrix}$$

$$(w \quad 1b) \underbrace{\begin{pmatrix} I & 0 \\ 0 & 0 \end{pmatrix}}_D \begin{pmatrix} w \\ b \end{pmatrix}$$

$$\begin{pmatrix} \epsilon_1 x_1 \dots & | & \epsilon_1 \\ \epsilon_2 x_2 \dots & | & \epsilon_2 \\ \vdots & & \vdots \\ \epsilon_N x_N \dots & | & \epsilon_N \end{pmatrix}$$

$$\begin{pmatrix} v \\ \vdots \\ b \end{pmatrix}$$

$$\geq \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$$

A

c

