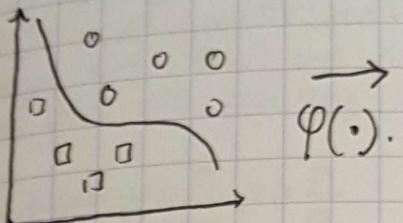


KERNELS (CSC411 NOTES VS SLIDES) (cannot read)

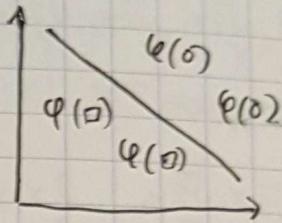
non-linear decision boundaries.

Feature mapping: Take input $x \rightarrow \phi(x)$.

INPUT SPACE:



FEATURE SPACE:



IF $x_i = [x^0, x^1, \dots, x^d]$, d high = HIGH COMP. COST.

QUADRATIC FEATURE BOUNDARY IDEA:

For $x = (1, x_1, x_2, \dots, x_d)$ we have:

$$\phi(x) = (1, \sqrt{2}x_1, \dots, \sqrt{2}x_d, \sqrt{2}x_1x_2, \sqrt{2}x_1x_3, \dots, \sqrt{2}x_1x_d, x_1^2, \dots, x_d^2)$$

Example: $d = 3$

$$x = (1, x_1, x_2, x_3)$$

$$\phi(x) = (1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_3, \sqrt{2}x_1x_2, \sqrt{2}x_1x_3, \sqrt{2}x_2x_3, x_1^2, x_2^2, x_3^2)$$

NOTE: dimension of $\phi(x)$ is thus order $O(d^2)$.

KERNELS

Recall: for $x = (1, x_1, x_2, x_3)$,

$$\phi(x) = (1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_3, \sqrt{2}x_1x_2, \sqrt{2}x_1x_3, \sqrt{2}x_2x_3, x_1^2, x_2^2, x_3^2)$$

Then: DEFINE K AS INNER PROD OVER $\phi(\cdot)$

So,

$$K(x, y) = \langle \phi(x), \phi(y) \rangle$$

$$K(x, y) = \langle \phi(x), \phi(y) \rangle$$

And, $\langle \phi(x), \phi(y) \rangle =$

$$\begin{aligned} & \langle (1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_3, \sqrt{2}x_1x_2, \sqrt{2}x_1x_3, \sqrt{2}x_2x_3, x_1^2, x_2^2, x_3^2, \\ & (1, \sqrt{2}y_1, \sqrt{2}y_2, \sqrt{2}y_3, \sqrt{2}y_1y_2, \sqrt{2}y_1y_3, \sqrt{2}y_2y_3, y_1^2, y_2^2, y_3^2, y_1^2y_2^2, y_1^2y_3^2, y_2^2y_3^2) \rangle \end{aligned}$$

The cross product:

$$\begin{aligned} 1 \cdot 1 &= 1, \quad \sqrt{2}x_1 \sqrt{2}y_1 = 2x_1y_1, \quad \sqrt{2}x_1x_2 \sqrt{2}y_1y_2 = 2x_1x_2y_1y_2 \\ \sqrt{2}x_2 \sqrt{2}y_2 &= 2x_2y_2, \quad \sqrt{2}x_1x_3 \sqrt{2}y_1y_3 = 2x_1x_3y_1y_3 \\ \sqrt{2}x_3 \sqrt{2}y_3 &= 2x_3y_3, \quad \sqrt{2}x_2x_3 \sqrt{2}y_2y_3 = 2x_2x_3y_2y_3 \end{aligned}$$

Result: $K(x, y) = (1 + \langle x, y \rangle)^2$.

computes K_{ij} , $O(d)$ time/memory, K is polynomial kernel.

(WHY IS THIS RELEVANT TO).

KERNELS: polynomial, gaussian, sigmoid

KERNEL PROPERTIES:

- 1) POS. CONST. $F \times N$. IS A KERNEL:
FOR $\alpha \geq 0$, $K'(\mathbf{x}_i, \mathbf{x}_j) = \alpha$.
- 2) POS. WEIGHTED LIN COMBS. OF KERNELS ARE KERNELS:
 $\forall i, \alpha_i \geq 0 : K'(\mathbf{x}_i, \mathbf{x}_j) = \sum_i \alpha_i K_i(\mathbf{x}_i, \mathbf{x}_j)$
- 3) PRODUCTS OF KERNELS ARE KERNELS
 $K'(\mathbf{x}_i, \mathbf{x}_j) = K_1(\mathbf{x}_i, \mathbf{x}_j) K_2(\mathbf{x}_i, \mathbf{x}_j)$

KERNEL LEGOS: combine to make complex feature maps!

Kernel feature
curves!

KERNELS

Geom. M.
For. Inv

OPTIMA

1) OP

2)

3)

Kernel feature spaces: become EXTREMELY complex, twisty,
curvy!

KERNELS AND SVM (CS 229 NOTES)

$$y_i = y_i \left(\left(\frac{w}{\|w\|} \right)^T x_i + \frac{b}{\|w\|} \right)$$

\geq margin.

Geom. Margin wrt (w, b) and x_i : $\hat{\gamma} = \min_i y_i$ over $i=1 \dots m$
for the m data pts resulting in $\gamma_1 \dots \gamma_m$ margins.

OPTIMAL MARGIN.

1) Optimization I:

$$\begin{aligned} & \text{maximize } \gamma \text{ wrt } \gamma, w, b \\ & \text{w/ constraints: } \|w\| = 1 \quad (\text{ensures func. marg. = geom. marg.}) \\ & \quad y_i(w^T x_i + b) \geq \gamma \quad \forall i = 1 \dots m. \end{aligned}$$

PROBLEM: $\|w\|=1$ is a non-convex constraint.

2) Optimization II:

$$\begin{aligned} & \text{maximize } \hat{\gamma}/\|w\| \text{ wrt } \hat{\gamma}, w, b \\ & \text{w/ constraints: } y_i(w^T x_i + b) \geq \hat{\gamma}, \quad \forall i = 1 \dots m. \end{aligned} \quad (\hat{\gamma} = \text{func. marg.})$$

PROBLEM: $\hat{\gamma}/\|w\|$ is a non-convex objective fn.

3) Optimization III:

Set $\hat{\gamma} = 1$. Then $\hat{\gamma}/\|w\| = 1/\|w\|$, which maximizing is
the same as minimizing $\|w\|^2$. Then:

$$\begin{aligned} & \text{minimize } \frac{1}{2} \|w\|^2 \\ & \text{w/ constraints: } y_i(w^T x_i + b) \geq 1 \quad \forall i = 1 \dots m. \end{aligned}$$

Solve w/ QP code or Lagrange duality.

LAGRANGE DUALITY

Consider minimize $f(w)$ wrt w
w/ constraint $h_i(w) = 0 \quad i=1 \dots m$.

LAGRANGIAN:

$$L(w, \beta) = f(w) + \beta_1 h_1(w) + \dots + \beta_m h_m(w)$$

β_i 's = Lagrange Multipliers.

Solve for w and β via: $\frac{\partial L}{\partial w_i} = 0 \quad \forall i = 1 \dots m, \quad \frac{\partial L}{\partial \beta_i} = 0, \quad \forall i = 1 \dots m$

PRIMAL:

$$\min f(w) \text{ wrt } w.$$

$$\text{s.t. } g_i(w) \leq 0 \quad i = 1 \dots n$$

$$h_i(w) = 0 \quad i = 1 \dots m.$$

GENERALIZED LAGRANGIAN:

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

$$\sum_{i=1}^m \beta_i h_i(w)$$

$$L(w, \alpha, \beta) = f(w) + \underbrace{\sum_{i=1}^n \alpha_i g_i(w)}_{\text{Primal}} + \dots + \alpha_n g_n(w) + \underbrace{\sum_{i=1}^m \beta_i h_i(w)}_{\text{Dual}} + \dots + \beta_m h_m(w).$$

Define $\Theta_P(w)$ as.

$$\Theta_P = \begin{cases} f(w) & \text{if } w \text{ satisfies primal constraints} \\ \infty & \text{else.} \end{cases}$$

whereas $\Theta_D(\alpha, \beta) = \max_w L(w, \alpha, \beta)$ wrt α, β , and $\alpha_i \geq 0$.

DUAL: $\Theta_D(\alpha, \beta) = \min_w L(w, \alpha, \beta)$

Then:

$$\alpha^* = \max_{\alpha \in \mathbb{R}^n} \min_w L(w, \alpha, \beta) \leq \underbrace{\min_w \max_{\alpha, \beta} L(w, \alpha, \beta)}_{\Theta_D(\alpha, \beta)} = p^*.$$

THUS: WE CAN SOLVE DUAL IN LIEU OF THE PRIMAL.

WE assume $d^* = p^*$.

Thus, maximizing dual \leftrightarrow minimizing primal.

thus, must exist w^*, α^*, b^* s.t. w^* soln of PRIMAL, α^*, b^* satisfy KKT condns:
 DUAL: $p^* = \partial^T = L(w^*, \alpha^*, b^*)$, and w^*, α^*, b^* satisfy KKT soln of
 $\alpha^* \geq 0, g_i(w^*) \leq 0, \underbrace{\alpha_i^* g_i(w^*)}_\text{if } \alpha_i^* > 0 = 0, \partial_{\beta} L = 0, \partial_{w_i} L = 0$.

OPTIMAL MARGIN CLASSIFIERS

recall $\min \frac{1}{2} \|w\|^2$ wrt w, b
 sub. $y_i(w^T x_i + b) \geq 1$.

Rewrite constraint as: $g_i(w) = -y_i(w^T x_i + b) + 1 \leq 0 \quad i=1 \dots n$

NOTE: $x_i \geq 0$ only for training examples s.t. their functional margin = 1, b/c they correspond to constraints $g_i(w) = -y_i(w^T x_i + b) = 0$.

They define SUPPORT VECTORS

Lagrangian for this optimization problem:

$$L(w, \alpha, b) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i (y_i(w^T x_i + b) - 1)$$

(No b_i b/c only inequality constraints on this optimal margin problem).

Taking derivative wrt w, b , we can get the dual problem wrt α .

KERNELS (CS 229)

Given "attributes" for each of x_i , we can use a REDUCE MAPPING $\phi(x)$ to get nonlinear features.

From before: if we have found the α_i 's, to make a prediction, we need to calculate:

$$w^T x + b = \left(\sum_{i=1}^m \alpha_i y_i x_i \right)^T x + b \\ = \boxed{\sum_{i=1}^m \alpha_i y_i \langle x_i, x \rangle + b}$$

For our new $\phi(x)$ features, we just replace each ~~x_i~~ , x_i , $\phi(x_i)$, $\phi(x)$. THUS,

$$w^T \phi(x) + b = \sum_{i=1}^m \alpha_i y_i \langle \phi(x_i), \phi(x) \rangle + b$$

DEFINE: KERNEL: $K(x, z) = \langle \phi(x), \phi(z) \rangle = \phi(x)^T \phi(z)$.

$\phi(x)$ might be hard to calculate, but $K(x, z)$ easy to calculate. e.g. $\phi(x)$ is EXTREMELY HIGH DIMENSION.

By calculating $K(x, z)$ efficiently, we can get an SVM to learn in ϕ w/o ever calculating $\phi(x)$ which would be hella inefficient.

SPANISH:

YO ~~USTE~~
to

THREE R

- OR
- ER
- IT

- AR E

hablo

- ER

com

- IR

vivo

TWO

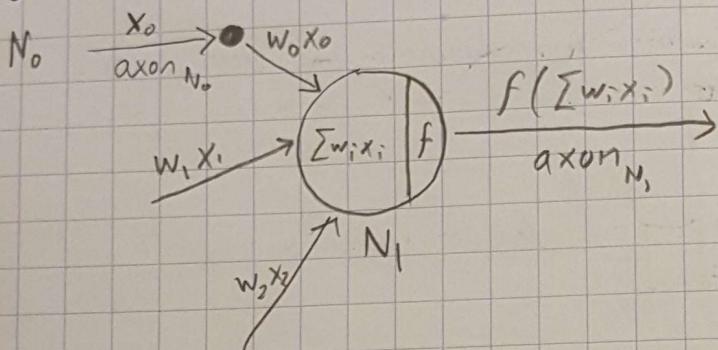
SU
C

- * MULTICLASS CLASSIFICATION. - REVIEW.
- * CH 4 → LINEAR MODELS FOR CLASSIFICATION.

NEURAL NETWORKS I

- + nonlinear discrimin. clfs using large collections of simple functions composed together.

Diagram of a neuron:



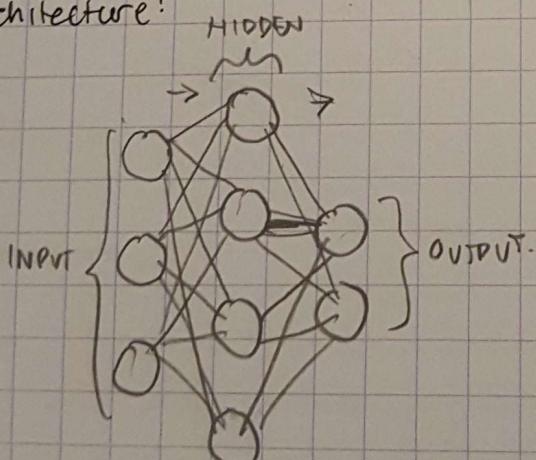
- + The dendrite inputs are "HIDDEN FEATURES"
- + Common activation functions:

+ sigmoid: $\sigma(z) = \frac{1}{1 + \exp(-z)}$

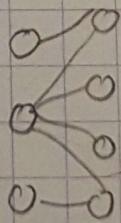
+ tanh: $\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$

+ ReLU: $\text{ReLU}(z) = \max(0, z)$

Architecture:



- + Here, 4 hidden-units
- + OUTPUT OF UNIT: LIN COMB of outputs from its inputs, with an activation function applied on it
- + NOTE:



is possible, aka we don't need all unit inputs from $i-1^{\text{th}}$ layer for all units in i^{th} layer.

+ NN LAYERS: we
+ REPRESENTATION
+ NN's are
+ Eg, $f_i(x)$
+ KEY: the
+ An NN w/
a universal
discriminat...
+ ↑ LAYER
run th...
PASSES:
+ FORWARD
+ BACKWARD
FORWARD
+ LEGI...

NN LAYERS: WE EXCLUDE the input layer. So the prev. arch. was
A 2-LAYER NN.

REPRESENTATIONAL POWER!

- + NN's are based on repeated, chained function composition.
- + E.g., $f_1(x)$, $f_2(x)$ and the output of N_3 is $f_3(w_1 f_1(x) + w_2 f_2(x))$, and f_3 may be the input for another neuron.
- + KEY: the weighting of the inputs (in the dendrites).
- + An NN w/ ONE HIDDEN LAYER can represent ANY FUNCTION as a universal approximator. Thus, good base for nonlinear discriminants.
- + ↑ LAYERS, ↑ HIDDEN UNITS = ↑ network capacity, but can run the risk of overfits.

PASSES:

- + FORWARD PASS: Make the prediction.
- + BACKWARD PASS: Compute the gradient

FORWARD PASS: COMPUTING A PREDICTION ON x .

- + LEGEND: j will index hidden units $h_j(x) = \text{out of } HU_j$
 k will index output units $O_k(x) = \text{out of out } u$.
 D is num inputs.
 v_{ji} is a weight from INPUT i to HU_j .
 w_{kj} is a weight from HU_j to OUT k .
- + OUTS: $h_j(x) = f(v_{j0} + \sum_{i=1}^D x_i v_{ji})$.
 $O_k(x) = g(w_{k0} + \sum_{j=1}^H h_j(x) w_{kj})$.

- + Example: $x = (x_1, x_2, x_3)$, $J = 4$, $O = 2$

Then,

$$h(x) = (h_1(x), h_2(x), h_3(x), h_4(x))$$
$$O(x) = (O_1(x), O_2(x))$$

$$h_j(x) = f(v_{j0} + x_1 v_{j1} + x_2 v_{j2} + x_3 v_{j3})$$
$$O_4(x) = g(w_{k0} + h_1(x) w_{k1} + h_2(x) w_{k2} + h_3(x) w_{k3})$$

Thus, O , h , x are vectors!

+ IF we don't use an activ. fxn, it's just linear classification b/c we are only looking at linear combinations. (More explicitly: linear regression).

+ SINGLE LAYER (NO HIDDEN) + SIGMOID = THIS IS JUST LOGISTIC REGRESSION

BACKWARD PASS: TRAINING THE NN VIA BACKPROP + ∇ s.

+ Usually to optimize a classifier, we find W weight vector such that this weight vector minimizes our chosen loss function.

+ BUT: hidden units = objective fn NO LONGER CONVEX. Empirically we know that many local minima are good though so we can still try SGD.

+ Enter BACKPROP. For each input x we need to do our classifying forward pass, THEN PROPAGATE GRADIENTS BACK and update W s via gradient descent.

BACKWARD PASS AND CHAIN RULE.

+ PROBLEM: when optimizing W , we need targets (aka test set). But hidden units don't give us targets, only output units do. Why?

+ Compute HOW FAST ERROR CHANGES AS ACTIVITY CHANGES.

+ We use ERROR DERIVATIVES WRT HIDDEN ACTIVITIES

+ We need to COMBINE THE EFFECTS each hidden activity has on the error.

+ Once we have error derivatives for all HAs, we can get error derivs. for the weights entering a HA.

↪ CHAIN RULE.

$$\partial \text{sig} = \sigma(z)(1-\sigma(z))$$

$$\partial \tanh = 1/\cosh^2(z)$$

$$\partial \text{ReLU} = 1, \text{if } z > 0, 0 \text{ if } z \leq 0.$$

SINGLE LAYER NET

ERROR GRADIENTS: $\frac{\partial E}{\partial w_{ki}} = \frac{\partial E}{\partial o_u} \frac{\partial o_u}{\partial z_k} \frac{\partial z_k}{\partial w_{ki}}$ where z_k is the combined INPUT to o_u before the activation func.

Some notation: $\delta_u^0 = \frac{\partial E}{\partial o_u}$
 $\delta_u^z = \delta_u^0 \cdot \frac{\partial o_u}{\partial z_k}$

and we know: $\frac{\partial o_u}{\partial z_k} z_k$ is just x_i , so $\frac{\partial o_u}{\partial z_k} F = S_u^z \cdot x_i$.

CONCRETE EXAMPLE:

- + Let E be MSE and g be LOGISTIC SIGMOID A.F.
- + Then...

$$\frac{\partial o_u}{\partial z_k} E = \frac{\partial o_u}{\partial z_k} \frac{1}{2} (o_u^{(n)} - t_u^{(n)})^2 = o_u^{(n)} - t_u^{(n)} \text{ for example (n),}$$

so $\delta_u^0 = o_u^{(n)} - t_u^{(n)}$

Next, we know $o_u^{(n)}$ is logistic sigmoid, so

$$o_u^{(n)} = g(z_u^{(n)}) = (1 + \exp(-z_u^{(n)}))^{-1}$$

$$\frac{\partial o_u}{\partial z_k} o_u = o_u^{(n)} (1 - o_u^{(n)})$$

$$\text{So, } S_u^z = \delta_u^0 \cdot \frac{\partial o_u}{\partial z_k} o_u = (o_u^{(n)} - t_u^{(n)}) o_u^{(n)} (1 - o_u^{(n)})$$

AND OUR FULL GRADIENT for $x^{(n)}$ is:

$$\frac{\partial o_u}{\partial w_{ki}} E = \sum_{n=1}^N (o_u^{(n)} - t_u^{(n)}) o_u^{(n)} (1 - o_u^{(n)}) x_i^{(n)}$$

AND we recall the GRAD UPDATE RULE:

$$w_{ki} \leftarrow w_{ki} - \eta \frac{\partial o_u}{\partial w_{ki}} E$$

EXAMPLE WITH A HIDDEN LAYER:

- + Legend:
 - let O_u be the outputs
 - let Z_u be the combined lincomb inputs to out u
 - let g be the output AF.
 - let w_{uj} be the weight of HU j on OUT u .

- let h_j be the HU output
- let u_j be the combined lincomb inputs to HU j .
- let f be the HU AF.
- let v_{ji} be the weight of INPUT i on HU j

- let X_i be INPUT i .

We can treat the $\partial_{w_{uj}} E$ as our single layer. But, instead of X_i , it's now h_j . Thus:

$$\begin{aligned}\partial_{w_{uj}} E &= \sum_{n=1}^N \underbrace{\partial_{O_u^{(n)}} E}_{\delta_u^{(n)}} \underbrace{\partial_{Z_u^{(n)}}}_{\partial_{w_{uj}}} \underbrace{\partial_{w_{uj}} Z_u^{(n)}}_{h_j^{(n)}} \\ &= \sum \delta_u^{(n)} h_j^{(n)}\end{aligned}$$

What about hidden weight gradients?

We use BACKPROP: FOR A SINGLE EXAMPLE $X^{(n)}$, calc $\partial_{w_{uj}} E$:

$$\begin{aligned}\partial_{h_j^{(n)}} E &= \sum_{k=1}^K \underbrace{\partial_{O_k^{(n)}} E}_{\delta_k^{(n)}} \underbrace{\partial_{Z_k^{(n)}}}_{\partial_{w_{uj}}} \underbrace{\partial_{w_{uj}} Z_k^{(n)}}_{w_{uj}} = \sum_{k=1}^K \delta_k^{(n)} w_{uj} \\ \text{call } \sum_{k=1}^K \delta_k^{(n)} w_{uj} &= \delta_j^{h,(n)}\end{aligned}$$

And Finally...

$$\begin{aligned}\partial_{v_{ji}} E &= \sum_{n=1}^N \underbrace{\partial_{h_j^{(n)}} E}_{\delta_j^{h,(n)}} \underbrace{\partial_{u_j^{(n)}}}_{\delta_j^{h,(n)}} \underbrace{\partial_{v_{ji}}}_{f'(u_j^{(n)})} \underbrace{u_j^{(n)}}_{x_i^{(n)}} = \sum_{n=1}^N \delta_j^{h,(n)} f'(u_j^{(n)}) x_i^{(n)}\end{aligned}$$

MULTIPLE HIDDEN LAYERS.

- + use the same idea as over there. For layer L , we need to backprop: compute $\partial_{y^L} E$ and we can use it to compute $\partial_{v_{ij}^L} E$ (weights at layer L) and $\partial_{h^{L-1}} E$.
- + Training NNs is HARD: long compute times, vanishing gradient problem (the gradient will be vanishingly small, preventing the weight from changing its value, which can even stop the NN from more changing).
- + ReLU is a solution to van-grad, and tanh/sigmoid can saturate causing vanishings. However, neurons can still "die" w/ ReLU...
- + OUTPUT g : either sigmoid/softmax, (classification), or NO output (regression).

+ WEIGHT INITIALIZATIONS:

- + CONSTANT \rightarrow all neurons will stay same
 - + STANDARD $\rightarrow w_{ij} \sim N(0, \sigma^2)$: $\sigma^2 \downarrow$, converge to 0, $\sigma^2 \uparrow$, diverge.
 - + XAVIER $\rightarrow \sigma^2 = 2/n_m + n_{out}$
 n_m, n_{out} # of units in prev layer and next layer
 - + HE $\rightarrow \sigma^2 = 2/n_m$.
Rec'd for ReLUs.
- Choosing σ^2
- + NEED MOMENTUM: $v_{t+1} = \beta v_t + \nabla L(w_t)$. $\beta = 0.9$.
 $x_{t+1} = \alpha x_t - \alpha v_{t+1}$.

NEURAL NETWORKS II.

= depend
INNOVUS c
re some
) = Γ_y
 $(x) = w$
 $D(x_d)$
gives
Good
VARI
Y CI

- + some reasons why computers SUCK at scene/obj recognition:
OCCLUSION, SCALE, DEFORM, CLUTTER, LIGHTING, VIEWPT, POSE,
HUGE IN-CLASS VARIATIONS
TONS OF CLASSES.

NEURAL NETS + IMAGES.

PREVENT OVERFITTING

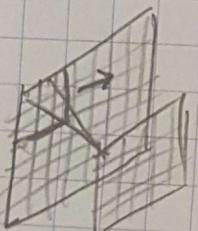
- + limit hidden unit #
- + limit weight size (weight decay) \rightarrow DROPOUT
- + Early stopping

- + LOCALLY-CONNECTED LAYERS: Images have many PIX, also \times is high dim. Fully-connected layers unfeasible. Good when input image is registered. (faces). Statistics similar at diff locs.
- + DIFFERENT FILTER FOR EACH (x, y) POSITION. 4M PARAMS ...

- + REPLICATED FEATURE APPROACH (monkey vision):

- + Many diff copies of same feature detector.
- + Use several DIFF FEAT TYPES each w/ own replicated pool of detectors.

- + CNN: connect each hidden unit to SMALL INPUT PATCH and SHARE WEIGHT ACROSS SPACE.



CNN hyperparams: DEPTH (# filters)

STRIDE (units apart a filter is applied)

SIZE (w x h) of filters

Need to pad output to prevent shrinkage.

- + POOLING: MAX or AVG. of features at specific locations.

1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

Max pool 2x2
Filter stride=2

6	8
3	4

- + FILTERBANK \rightarrow CUBE. Each slice is the output of convolving img w/ one filter. APPLY activation fn on each hidden unit.
- + MAX POOLING (robustness) (INVARIANCE IN SMALL POSITION SHIFTS).

- + # HIDDEN US smaller than input or output units = not generalized, lost info. keep them bigger! (can be used 4 autoencoders).

PRINCIPAL COMPONENTS ANALYSIS.

SOME GOALS OF UNSUPERVISED LEARNING

- + reducing data dimensionality (while retaining its key nature)
- + finding CLUSTERS (data w/ similarities, forming classes)
- + modelling data density
- + finding hidden causes.

1. DIM RED \rightarrow rep each input case w/ small number of vars

2. CLUSTER \rightarrow rep each input case w/ PROTOTYPE EXAMPLE

3. DENSITY EST \rightarrow estim PROB DIST over the data

+ AIM OF PCA: find small # OF "DIRECTIONS" (axes) in input space that EXPLAIN INPUT DATA VARIANCE.

We then RE-PROJECT DATA ONLY ON THOSE AXES.

+ Assume data is continuous

+ CAN REDUCE OVERRIFTS

+ Good for visualization, preprocessing, lossy compression, better generalization.

let $K \ll D$, and $\mathbf{z} \in \mathbb{R}^K$, so goal: $\mathbf{x} \approx \mathbf{U}\mathbf{z} + \mathbf{a}$ and we project \mathbf{x} to lower dimensions. \mathbf{U} is a $D \times K$ MATRIX.

+ SEARCH for [ORTHOGEN DIRS IN SPACE w/ HIGHEST VAR]

- + CENTER the data. (FIND \bar{x} , subtract $x - \bar{x}$)
- + CALCULATE EMPIRICAL COVARIANCE MATRIX:

$$\Sigma = \frac{1}{N} X^T X \quad \text{Note: } X^T X = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \\ \vdots & \vdots \\ A_{n1} & A_{n2} \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ \vdots & \vdots & \vdots \\ A_{n1} & A_{n2} & A_{n3} \end{bmatrix}$$

And: Result: 3×3 matrix

w/ elements: $A_{11}A_{11} + A_{12}A_{21}$ $A_{11}A_{22} + A_{12}A_{22}$, etc.

- + Look for a direction/axis w that MAXIMIZES PROJECTION VARIANCE $y_i = w^T x_i$. Normalize: $\|w\|=1$.

$$+ \text{Var}(y) = w^T \Sigma w = \sum_i \frac{1}{N} (w^T x_i)^2 = \frac{1}{N} \sum_i w_i^2 x_i^T w$$

- + So, we want to maximize the variance and find the direction that maximizes variance. Thus:

$$w^* = \underset{\|w\|=1}{\operatorname{argmax}} w^T \Sigma w.$$

HOW TO MAXIMIZE AND FIND w^* :

- 1) Σ has an EIGENDECOMP. w/ ORTHONORMAL BASIS VECTORS
- 2) Call them v_1, \dots, v_d , and $\lambda_1 \geq \dots \geq \lambda_d \geq 0$.
- 3) Write w in terms of those bases:

$$w = \sum_i a_i v_i \quad \text{and} \quad \sum_i a_i^2 = 1.$$

- 4) This transforms our objective fn. Previously we wanted to maximize

$$w^T \Sigma w.$$

But b/c of eigendecomp now we need to maximize

$$\sum_{i: a_i^2=1} a_i^2 \lambda_i$$

TOP k DIMENSIONS ARE TOP k EIGENVECTORS

Another way to see:

+ $\Sigma = U\Lambda U^T$ and $U^T U = U U^T = I$
+ Take $Z^i = U_u^T x^i$ Take first k eigenvectors

+ therefore, $Z = U_{1:k}^T (x - m)$ (where m is mean)

+ To project, we calculate empirical cov matrix and
find its eigenvectors. That's PCA in a nutshell.

+ Approximate reconstruction of $\tilde{x} = U_{1:u} Z + m$

FINDING w_i s w/ Lagrange multipliers:

1) FINDING w_i :

GOAL: Maximize $w_i^T \Sigma w_i$ and $\|w_i\| = 1$

Let α be a Lagrange multiplier. And $w_i^T w_i = 1 \rightarrow 0 = 1 - w_i^T w_i$

Then:

$$w_i^T \Sigma w_i + \alpha(1 - w_i^T w_i)$$

Take derivative wrt w_i :

$$\sum w_i + \alpha w_i = 0$$

$\sum w_i = \alpha w_i$ so w_i is eigenvector of Σ w/
eigval α .

This also shows w_i must be the eigenvector w/ max eigenval
of max value b/c this maximizes objective

2) FINDING w_2

GOAL: Maximize $w_2^T \Sigma w_2$ w/ $\|w_2\| = 1$ and $w_2^T w_1 = 0$.

LAGRANGE:

$$w_2^T \Sigma w_2 + \alpha_2(1 - w_2^T w_2) - \beta w_2^T w_1$$

Derivative wrt w_2 :

$$\sum w_2 - \alpha_2 w_2 - \beta w_1 = 0$$

$$w_1^T \Sigma w_2 - \alpha_2 w_1^T w_2 - \beta w_1^T w_1 = 0$$

$$\alpha_2 w_1^T w_2 - \alpha_2 w_1^T w_2 - \beta w_1^T w_1 = 0$$

$$\alpha_2 0 - \alpha_2 0 - \beta = 0$$

$$\beta = 0$$

Therefore, $w_2^T \Sigma w_2 + \alpha_2(1 - w_2^T w_2)$ results in

$$\sum w_2 = \alpha_2 w_2$$

y_i = dependent variable
continuous or -

t-SNE + STOCHASTIC NEIGHBOUR EMBEDDING

+ PCA tries to find GLOBAL STRUCTURE so its bad for fine details, can lead to local inconsistencies (far pt = nearest neighbour)

+ t-SNE FINDS LOCAL STRUCTURE.

But is almost only used for visualization as a result b/c can't easily embed new pts.

SNE BASIC IDEAS

- + Encode the HDim data as distro, do a "random walk thru data, ↑ prob to jump to closer pts, and find Low Dim pts w/ sim neighbourhoods"
- + Distribution distance: KL distance
- + E.g. Gaussian distro. x_i pt, choosing x_j pt will be prob less as farther away

$$P_{j|i} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2 / 2\sigma_i^2)}$$

prob of choosing pt x_j from x_i , $P_{j|i} = 0$.

σ_i → sets size of neighbourhood.

σ_i → diff for each x_i

+ Find distro, symmetrize over pairs: $P_{ij} = 1/2N(P_{ilj} + P_{jli})$

+ PERPLEXITY: $\text{perp}(P_{j|i}) = 2^{H(P_{j|i})}$

$$H(P_{j|i}) = - \sum_k P_{j|ik} \log(P_{j|ik})$$

P uniform over k elems = perp is k.

Low perp = $\downarrow \sigma^2$

High perp = $\uparrow \sigma^2$

Perplexity: How well prob distro predicts sample
5-50 good.

+ When we have distro P w/ P_{ij} we can now define Q of Q_{ij}

$$Q_{ij} = \frac{\exp(-\|y^i - y^j\|^2)}{\sum_k \sum_{l \neq k} \exp(-\|y^i - y^k\|^2)}$$

Optimizing Q close to P means minimizing KL-divergence.
 y^1, y^2, \dots, y^N are embeddings we are optimizing

$$KL(Q||P) = \sum_{ij} Q_{ij} \log \left(\frac{Q_{ij}}{P_{ij}} \right)$$

distribution distance.

$KL(Q||P)$ is the "penalty" for using the wrong distro of expecting P but getting Q . $KL(Q||P) \geq 0$, $=0$ when $Q=P$.

$KL(Q||P)$ is CONVEX. If $P_{ij}=0$ but $Q_{ij}>0$, then

$$KL(Q||P) = \infty \text{ (VERY BAD)}$$

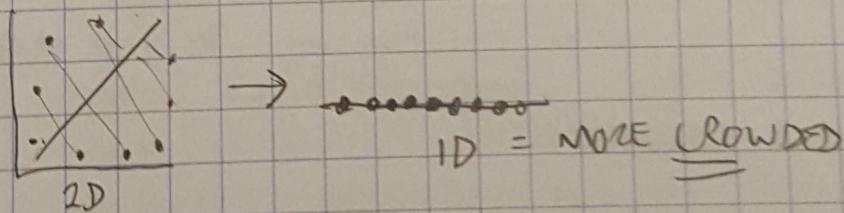
+ SNE : We have P .

looking for $y^1 \dots y^N \in \mathbb{R}^D$ s.t. our inferred distro Q will MINIMIZE $L(Q) = KL(P||Q)$.

$$\text{Can show } \partial_{y^i} L = \sum_j (P_{ij} - Q_{ij}) (y^i - y^j)$$

Not convex.

PROBLEM: CROWDING : As we embed in lower dimensions, less space, shorter distances!



t-SNE: change SNE Gaussian of Q to a t-DISTRIBUTION (heavy-tailed Gaussian)

Slower change in $Q \rightarrow$ more "room"

$$Q_{ij} = (1 + \|y_i - y_j\|^2)^{-1} \text{ instead of } \exp(-\|y_i - y_j\|^2).$$

w/ t-Distro.

$$\partial_{y^i} L = \text{same but times } (1 + \|y_i - y_j\|^2)^{-1}.$$

and y_i = dependent
CONTINUOUS

CLUSTERING.

- + Grouping N examples into K clusters
- + We initially assume the data belongs to K classes.
- + We assume data pts from the same class are similar (close Euclidean dist)
- + K-MEANS: start randomly w/ K centers and try using
 - + Initialization
 - + ASSIGNMENT STEP : Each data pt assigned to closest mean
 - + REFITTING STEP : Adjust the means.

$$\mu_k = \frac{\sum r_k^{(n)} x^{(n)}}{\sum r_k^{(n)}}$$

$r_k^{(n)}$ = responsibilities

ENSEMBLES. : A return to supervised learning.

- + ENSEMBLE OF PREDICTORS: a set of predictors whose INDIVIDUAL DECISIONS combined to CLASSIFY NEW EXAMPLES
 1. Generate MULTIPLE CLASSIFIERS
 2. Each votes on TEST INSTANCE.
 3. Take MAJORITY/AVERAGE AS PREDICTION

- + BAGGING: train SEPARATE MODELS on OVERLAPPING TRAIN SET. Average predictions

- + BOOSTING: Sequential, ITERATIVE re-weight of training examples so current classifier can focus on HARD EXAMPLES

BIAS-VARIANCE DECOMPOSITION.

$$\text{Error} = (\text{bias})^2 + \text{variance} + \text{noise}$$

$$\mathbb{E}_t[(t - y(x; D))^2] = \mathbb{E}_t[(t - h^*(x) + h^*(x) - y(x; D))^2] =$$

$$\underbrace{\mathbb{E}_t[(t - h^*(x))^2]}_{\text{noise.}} - 2\mathbb{E}_t[(t - h^*(x))(\overbrace{h^*(x) - y(x; D)}^{\text{bias}})] + \mathbb{E}_t[(h^*(x) - y(x; D))^2]$$

$$\mathbb{E}_D[(h^*(x) - y(x; D))^2] = \mathbb{E}[(h^*(x) - \mathbb{E}[h^*(x; D)])^2] + \underbrace{\mathbb{E}_{D, x}[(y(x; D) - \mathbb{E}[y(x; D)])^2]}_{\text{variance}}$$

+ BAGGING:
+ Random
+ OOB ESTIMATES

+ EVEN IF A WINNER DOES NOT THIS). OR

BOOSTING:
+ BIAS REDUCES
+ ADABOOST
doing error

+ ADA

+ ALC

- + **BAGGING**: Sampling w/ replacement \rightarrow used for data paucity.
+ Random Forest. Goal = reduce variance, reduce overfitting.
- + **OOB ESTIMATION**: predict each training example using
ALL TREES that DID NOT CONTAIN THIS
in training data.

- + Even if a single model is great, used in most competition
- winners -
- + does not reduce bias. Need more randomness (Random Forest does this). OOB reduces the need of cross-validation.

BOOSTING

- + BIAS reduction method in contrast with BAGGING which reduces variance
- + **ADABOOST** (adaptive boosting): Take a "weak" classifier doing $< e$ better than chance, "BOOST" to get low training error

CLASSIFIER: $H(x) = \text{sign}(\sum_i \alpha_i h(x_i))$

where h_i are weak classifiers

SUM OF LINEAR CLASSIFIERS NOT LINEAR

- + ADABOOST IDEA: Each iteration, reweigh training sample.
(\uparrow weight to points wrongly classified)
Train new weak classifier, minimizing weighted acc.

- + ALGORITHM:
 - 1) INPUT $\{(x^n, t^n)\}_{n=1}^N$, example weights $D_m^n(x) = 1/N$
 - 2) $m = 1 \rightarrow M$
 - a) $h_m(x) = \text{weak learn } x, t, w$
 - b) Fit $h_m(x)$ by minimizing

$$J_m = \sum_{n=1}^N D_m^n [h_m(x^n) \neq t^n]$$

- c) weighted error rate, classifier coefficient:

$$\epsilon_m = P_m / \sum D_m^n$$

$$\alpha_m = \frac{1}{2} \log \frac{1-\epsilon_m}{\epsilon_m}$$

- d) UPDATE: $D_{m+1}^n = D_m^n \exp(-t^n \alpha_m h_m(x^n))$