



## Lecture 22

- high-dimensional spaces
- Random projections
- Pseudoinverse

Consider random pt.  $\vec{p} \sim N(\mathbf{0}, I) \in \mathbb{R}^d$

What is the distribution of its length?

- You'd think each component is close to 0  
 $\hookrightarrow$  this is wrong in higher dimensions

turns out majority of pts. are at approximately the same distance from the mean  
 $\hookrightarrow$  they lie on a thin shell

Squared distance:  $\|\vec{p}\|_2^2 = p_1^2 + \dots + p_d^2$

each  $p_i$  sampled independently from a univariate standard normal distr.

Sum of each component comes from a **chi-squared distribution**:

$$p_i \sim N(0, 1) \quad p_i^2 \sim \chi^2(1) \quad E[p_i^2] = 1 \quad \text{Var}(p_i^2) = 2$$

$$\Rightarrow E[\|\vec{p}\|^2] = d \cdot E[p_i^2] = d$$

$$\text{Var}(\|\vec{p}\|^2) = d \text{Var}(p_i^2) = 2d$$

$$\sigma(\|\vec{p}\|^2) = \sqrt{2d}$$

large  $d \rightarrow \|\vec{p}\|$  (concentrated in thin shell w/ radius  $\approx \sqrt{d}$ ) and a thickness  $\propto (2d)^{1/4}$

mean  $E[\|\vec{p}\|]$  not exactly  $\sqrt{d}$ , but close since std deviation  $\sqrt{2d}$  much smaller than  $d$   
 $\hookrightarrow$  same w/  $\sigma(\|\vec{p}\|)$

What about a uniform distribution?

- consider 2 spheres of radii  $r$  and  $r-\varepsilon$

- volume of outer ball proportional to  $r^d$

- vol. of inner proportional to  $(r-\varepsilon)^d$

$$\Rightarrow \frac{V_1}{V_0} = \frac{(r-\varepsilon)^d}{r^d} = \left(1 - \frac{\varepsilon}{r}\right)^d \approx e^{-\frac{\varepsilon d}{r}}$$

ratio small for a large  $d \Rightarrow$  most of volume comes from very thin shell of outer sphere  
e.g.  $\frac{\varepsilon}{r} = 0.1$  &  $d=100 \Rightarrow V_1$  has  $0.9^{100} = 0.0027\%$  of volume

$\rightarrow$  random points from uniform & Gaussian distributions in high dimensions almost always fall in some outer shell

$\Rightarrow$  in high dimensions, nearest neighbor & 1000-nearest neighbors don't differ by much

$\Rightarrow$  k-nearest neighbors & k-means clustering are less effective for large  $d$

## Angles Between Random Vectors

What is the angle  $\theta$  between  $\vec{p} \sim N(0, I_d)$  & an arbitrary  $\vec{q} \in \mathbb{R}^d$ ?

W/o loss of generality, let  $\vec{q} = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$  \* q's value doesn't matter since p is random in all directions

$$\cos \theta = \frac{\vec{p}^T \vec{q}}{\|\vec{p}\| \|\vec{q}\|} = \frac{p_1}{\|\vec{p}\|} \xrightarrow{\mu=0, \sigma=1} E[\cos \theta] = 0, \sigma(\cos \theta) \approx 1/\sqrt{d}$$

$\rightarrow$  large d  $\Rightarrow \cos \theta \approx 0$ ,  $\cos \theta \approx 0 \Rightarrow \theta \approx 90^\circ$

$\rightarrow$  Vectors close to being orthogonal

## Random Projection

- alternative to PCA as preprocessing for clustering, classification, and/or regression
- approximately preserves distances between points

Procedure: project onto random Subspace (instead of PCA Subspace)

- Sometimes preserves distance better than PCA
- Works best when you project high-dim. Space to medium-dim Subspace
- Similar distance  $\Rightarrow$  k-means & nearest Neighbour perform similarly but quicker

• Pick Small  $\epsilon$ , Small  $\delta$ , random Subspace  $S \subset \mathbb{R}^d$  of dimension k  
 $k = \left\lceil \frac{2 \ln(1/\delta)}{\epsilon^2 - \epsilon^3/3} \right\rceil$  doesnt depend on d

round up  $\rightarrow$  for any pt.  $\vec{q}$ , let  $\vec{\tilde{q}}$  be orthogonal projection onto S, multiplied by  $\sqrt{\frac{d}{k}}$

**Johnson-Lindenstrauss Lemma:** for any  $\vec{q}, \vec{v} \in \mathbb{R}^d$  distance between original pts.

$$(1-\epsilon) \|\vec{q} - \vec{v}\|^2 \leq \|\vec{\tilde{q}} - \vec{\tilde{v}}\|^2 \leq (1+\epsilon) \|\vec{q} - \vec{v}\|^2 \quad w/ \text{prob. } \geq 1 - 2e^{-\delta} \quad (\text{since this process is random})$$

$$\epsilon \text{ typically } \in [0.02, 0.5] \quad \delta \in [\frac{1}{n^3}, 0.05]$$

w/ these values, squared distance after projection may change by 2% to 50%  
experiment w/ k to find best speed-accuracy trade off

Small  $\delta, \delta \leq 1/n^2 \Rightarrow$  inter-sample-pt. distances small

$\hookrightarrow$  Subspace of dim.  $\Theta(\log n)$

- reducing  $\delta$  doesn't cost much, reducing  $\epsilon$  costs more

- w/ random projections, you can bring 1,000,000  $\rightarrow$  1,000 w/ 6% error

## Why does this Work?

Random projection of  $\vec{q} \cdot \vec{u}$  is like taking a random vector & selecting k components

↳ Mean of Squares of the k components approximates mean for whole population

How do we get a uniformly distributed random proj. direction?

- choose each comp. from univariate gaussian & normalize

How do we get a random subspace?

1) Choose k random directions

2) Use GS to make them mutually orthogonal

## Pseudoinverse, SVD

Back to supervised learning!

Suppose  $D \in \mathbb{R}^{n \times d}$  is diagonal

then we find  $D^+$  by transposing  $D$  & replacing every nonzero value w/ its reciprocal

$$\rightarrow D = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & \frac{1}{3} \\ 0 & 0 & 0 \end{bmatrix} \quad D^+ = \begin{bmatrix} \frac{1}{2} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad DD^+ = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad D^+D = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

→ If  $D$  has  $n$  nonzero diagonal entries, then  $D^+$  is its inverse,  $DD^+ = D^+D = I$

In general,  $DD^+$  &  $D^+D$  are always diagonal matrices w/ 1's & 0's

-  $DD^+D = D$ ,  $D^+DD^+$  and  $D^2D^+ = D$

Now, Pseudoinverse of arbitrary  $X \in \mathbb{R}^{n \times d}$  w/ SVD  $X = UDV^T$ ,  $\text{rk}(X) = \text{rk}(D)$

then  $X^+ = VD^+U^T$

Properties: 1)  $XX^+ = U(DD^+)U^T$  which is symm. PSD

2)  $X^+X = V(D^+D)V^T$  also symm. PSD

3) Same rank:  $D$ ,  $D^+$ ,  $DD^+$ ,  $D^+D$ ,  $X$ ,  $X^+$ ,  $XX^+$ ,  $X^+X$

4) if  $\text{rk}(X) = n$ , then  $XX^+ = I_{n \times n}$ ,  $X^+$  is right pseudoinverse

5) if  $\text{rk}(X) = d$ ,  $X^+X = I_{d \times d}$ ,  $X^+$  is left inverse

6)  $XX^+X = X \quad 7) X^+X^+X = X^+$

- \* Pseudoinverse always gives good solution to LS linear regression even when  $X^T X$  is singular
- Proof: Show  $\vec{U} = X^+ \vec{y}$  is a solution to  $X^T X \vec{U} = X^T \vec{y}$
- if normal equations have multiple solutions, then  $\vec{U} = X^+ \vec{y}$  is the min. norm sol.
- $X^+$  helpful when  $X^T X$  is singular since n2d pts. lie on subspace of feature space
- if  $X$  has small singular value, then its reciprocal is large, large effect on  $\vec{U}$ 
  - if  $\sigma_i = 0$ , that value has no effect on  $\vec{U}$
- if we have a very small  $\sigma_j$ , should we pretend it's 0?
  - ridge regression kinda does this